

# Probabilistic disease progression modeling to characterize diagnostic uncertainty: Application to staging and prediction in Alzheimer's disease

Marco Lorenzi<sup>a,b,\*</sup>, Maurizio Filippone<sup>c</sup>, Giovanni B. Frisoni<sup>d,e</sup>, Daniel C. Alexander<sup>f</sup>, Sebastien Ourselin<sup>b</sup>, for the Alzheimer's Disease Neuroimaging Initiative<sup>1</sup>

<sup>a</sup> Asclepios Research Project, Université Côte d'Azur, Inria, France

<sup>b</sup> Translational Imaging Group, Centre for Medical Image Computing, University College London, UK

<sup>c</sup> EURECOM, France

<sup>d</sup> Geneva Neuroscience Center, University Hospitals and University of Geneva, Switzerland

<sup>e</sup> IRCCS Fatebenefratelli, Brescia, Italy

<sup>f</sup> POND Group, Centre for Medical Image Computing, University College London, UK

## ARTICLE INFO

### Keywords:

Alzheimer's disease  
Diagnosis  
Disease progression modeling  
Gaussian process  
Clinical trials

## ABSTRACT

Disease progression modeling (DPM) of Alzheimer's disease (AD) aims at revealing long term pathological trajectories from short term clinical data. Along with the ability of providing a data-driven description of the natural evolution of the pathology, DPM has the potential of representing a valuable clinical instrument for automatic diagnosis, by explicitly describing the biomarker transition from normal to pathological stages along the disease time axis. In this work we reformulated DPM within a probabilistic setting to quantify the diagnostic uncertainty of individual disease severity in an hypothetical clinical scenario, with respect to missing measurements, biomarkers, and follow-up information. We show that the staging provided by the model on 582 amyloid positive testing individuals has high face validity with respect to the clinical diagnosis. Using follow-up measurements largely reduces the prediction uncertainties, while the transition from normal to pathological stages is mostly associated with the increase of brain hypo-metabolism, temporal atrophy, and worsening of clinical scores. The proposed formulation of DPM provides a statistical reference for the accurate probabilistic assessment of the pathological stage of de-novo individuals, and represents a valuable instrument for quantifying the variability and the diagnostic value of biomarkers across disease stages.

## Introduction

Neurodegenerative disorders (NDDs), such as Alzheimer's disease (AD), are characterized by the progressive pathological alteration of the brain's biochemical processes and morphology, and ultimately lead to the irreversible impairment of cognitive functions (Brookmeyer et al., 2007). The correct understanding of the relationship between the different pathological features is of paramount importance for improving the identification of pathological changes in patients, and for better treatment (Jack et al., 2010).

To this end, ongoing research efforts aim at developing precise models allowing optimal sets of measurements (and combinations of

them) to uniquely identify pathological traits in patients. This problem requires the definition of optimal ways to integrate and jointly analyze the heterogeneous multi-modal information available to clinicians (Young et al., 2013; Mwangi et al., 2014; Lorenzi et al., 2016). By consistently analyzing multiple biomarkers that to date have mostly been considered separately, we aim at providing a richer description of the pathological mechanisms and a better understanding of individual disease progressions.

Disease progression modeling (DPM) is a relatively new research direction for the study of NDD data (Fonteijn et al., 2011; Jedynak et al., 2012; Donohue et al., 2014; Younes et al., 2014; Bilgel et al., 2015; Schiratti et al., 2015; Guerrero et al., 2016; Marinescu et al., 2017). The

\* Corresponding author. 2004 Route des Lucioles, Inria Sophia Antipolis, 06902 Valbonne, France.

E-mail addresses: [marco.lorenzi@inria.fr](mailto:marco.lorenzi@inria.fr) (M. Lorenzi), [maurizio.filippone@eurecom.fr](mailto:maurizio.filippone@eurecom.fr) (M. Filippone), [Giovanni.Frisoni@unige.ch](mailto:Giovanni.Frisoni@unige.ch) (G.B. Frisoni), [d.alexander@ucl.ac.uk](mailto:d.alexander@ucl.ac.uk) (D.C. Alexander), [s.ourselin@ucl.ac.uk](mailto:s.ourselin@ucl.ac.uk) (S. Ourselin).

<sup>1</sup> Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: [http://adni.loni.usc.edu/wp-content/uploads/how\\_to\\_apply/ADNI\\_Acknowledgement\\_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf).

main goal of DPM consists in revealing the natural history of a disorder from collections of imaging and clinical data by: 1) *quantifying* the dynamics of NDDs along with the related temporal relationship between different biomarkers, and 2) *staging* patients based on individual observations for diagnostic and interventional purposes. Therefore, this research domain is closely related to the exploitation of advanced statistical/machine-learning approaches for the joint modeling of the heterogeneous and information available to clinicians: imaging, biochemical, and clinical biomarkers. Differently from the several predictive machine-learning approaches proposed in the past in NDD research, disease progression models aim at explicitly estimating the temporal progression of the biomarkers from normal to pathological stages, to provide a better interpretation and understanding of the natural evolution of the pathology. For this reason it represents a very appealing modeling approach in clinical settings.

The main challenge addressed by DPM consists in the general lack of well-defined temporal reference in longitudinal clinical dataset of NDDs. Indeed, age or visit date information are biased time references for the individual longitudinal measurements, since the onset of the pathology may vary across individuals according to genetic and environmental factors (Yang et al., 2011). This is a very specific methodological issue requiring the extension and generalization of the analysis approaches classically used in time-series analysis.

To tackle this problem, it is usually assumed that individual biomarkers are measured relatively to an underlying disease trajectory defined with respect to an absolute time axis describing the natural history of the pathology (Jedynak et al., 2012). Each individual is thus characterized by a specific observation time that needs to be estimated in order to assess the individual pathological stage. According to this statistical setting, we therefore aim at estimating a *group-wise* disease model defined with respect to an absolute time scale, along with *individual* time re-parameterisation relative to the group-wise progression. This modeling paradigm has been implemented in a number of approaches proposed in the recent years, either by assuming continuous temporal trajectories of the biomarkers (Jedynak et al., 2012; Donohue et al., 2014; Younes et al., 2014; Bilgel et al., 2015; Schiratti et al., 2015; Guerrero et al., 2016; Marinescu et al., 2017), or by modeling the disease progression as a sequence of discrete events (Fontejn et al., 2011; Young et al., 2014).

For example, in (Donohue et al., 2014) the authors proposed to model the temporal biomarker trajectories through random effect regression, building on the theory of self-modeling regression (Kneip and Gasser, 1988), while the authors of (Schiratti et al., 2015) re-frame the random effect regression model in a geometrical setting, based on the assumption of a logistic curve shape for the average biomarker trajectories.

Continuous progression models have been recently extended to the modeling of brain images based on the time-reparameterization of voxel/mesh-based measures (Younes et al., 2014; Bilgel et al., 2015; Marinescu et al., 2017).

The use of disease progression models for diagnostic purposes is instead less investigated. Predictive models of patient staging were proposed within the setting of the Event Based Model (Fontejn et al., 2011), or still through random effect modeling (Guerrero et al., 2016). However, the Event Based Model relies on the coarse binary discretization of the biomarker changes, and does not account for longitudinal observations, while the predictive models proposed in (Guerrero et al., 2016) and (Schmidt-Richberg et al., 2015) require cohorts with known disease onset, and therefore lack flexibility while being prone to bias due to mis-diagnosis and uncertainty of the conversion time.

Furthermore, these methods are generally not formulated in a probabilistic setting, which makes it difficult to account for uncertainties in biomarker progressions and diagnostic predictions. Indeed, the quantification of the variability associated with the biomarkers trajectories, as well as the assessment of the diagnostic uncertainty in *de-novo patients*, are crucial requirements for decision making in clinical practice (Shinkins and Perera, 2013).

Nonetheless, the ensemble of this research offers a sight of the potential of these approaches in representing a novel and powerful diagnostic instrument: in this study we thus aim at assessing the ability of DPM in providing a statistical reference for the transition from normal to pathological stages, for probabilistic diagnosis in the clinical scenario. To this end, we reformulate classical DPM within a Bayesian setting in order to allow the probabilistic estimate of the biomarker trajectories and the quantification of the uncertainty of predictions of the individual pathological stage. The resulting probabilistic framework is exploited in an hypothetical clinical scenario, for the estimation of the pathological stage in a *de-novo* cohort of testing individuals, by assessing the influence of missing observations, biomarkers, and follow-up information.

The manuscript is structured as follows. Section 2.1 formulates DPM based on Bayesian Gaussian Process regression (Rasmussen, 2006), while Section 2.2 illustrates the validation of our model on clinical and multivariate imaging measurements from a cohort of 782 amyloid positive individuals extracted from the ADNI database.

## Methods

### Statistical setting

This section highlights the statistical framework employed in this study, based on the reformulation of self-modeling regression with a Bayesian setting. This achieved by 1) defining a random effect Gaussian process regression model to account for individual correlated time series (section 2.1.1); 2) modeling individual time transformations encoding the information on the latent pathological stage (section 2.1.2); and 3) introducing a monotonicity information in order to impose a regular behaviour on the biomarkers trajectories (section 2.1.3). We finally illustrate in section 2.1.4 how the proposed framework leads to a probabilistic model of disease staging in *de-novo* individuals, naturally accounting for missing information. Further details on model specification and inference are provided in the Supplementary Section Appendix A.1, while the experimental validation on synthetic data is reported in Supplementary Section Appendix A.2.

### Gaussian process-based random effect modeling of longitudinal progressions

In what follows, longitudinal measurements of  $N_b$  biomarkers  $\{b_1, \dots, b_{N_b}\}$  over time are given for  $N$  individuals.

We represent the longitudinal biomarker's measures associated with each individual  $j$  as a multidimensional array  $(\mathbf{y}^j(t_1), \mathbf{y}^j(t_2), \dots, \mathbf{y}^j(t_{k_j}))^\top$  sampled at  $k^j$  multiple time points  $\mathbf{t} = \{t_1, t_2, \dots, t_{k^j}\}$ . Although different biomarkers may be in reality sampled at different time-points, for the sake of notation simplicity in what follows we will assume, without loss of generality, that the sampling time is common among them. The observations for individual  $j$  at a single time point  $t$  are thus a random sample from the following generative model:

$$\mathbf{y}^j(t) = \left( y_{b_1}^j(t), y_{b_2}^j(t), \dots, y_{b_{N_b}}^j(t) \right)^\top \quad (1)$$

$$= \mathbf{f}(t) + \boldsymbol{\nu}^j(t) + \boldsymbol{\varepsilon}, \quad (2)$$

where  $\mathbf{f}(t) = (f_{b_1}(t), f_{b_2}(t), \dots, f_{b_{N_b}}(t))^\top$  is the fixed effect function modeling the biomarker's longitudinal evolution,  $\boldsymbol{\nu}^j(t) = (\nu_{b_1}^j(t), \nu_{b_2}^j(t), \dots, \nu_{b_{N_b}}^j(t))^\top$  is the individual random effect, and  $\boldsymbol{\varepsilon} = (\varepsilon_{b_1}, \varepsilon_{b_2}, \dots, \varepsilon_{b_{N_b}})^\top$  is time-independent observational noise. The group-wise evolution is modelled as a GP,  $\mathbf{f} \sim \mathcal{GP}(0, \Sigma_G)$ , the individual random effects are assumed to be correlated perturbations  $\boldsymbol{\nu}^j \sim \mathcal{N}(0, \Sigma_S)$ , while the observational noise is assumed to be a Gaussian heteroskedastic term  $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma_\varepsilon)$ , where  $\Sigma_\varepsilon$  is a diagonal matrix  $\text{diag}[\sigma_{b_1}^2, \sigma_{b_2}^2, \dots, \sigma_{b_{N_b}}^2]$ .

*Fixed effect process.* The covariance function  $\Sigma_G$  describes the biomarkers

temporal variability, and is represented as a block-diagonal matrix

$$\Sigma_G(\mathbf{f}, \mathbf{f}) = \text{diag} \left[ \sum_{b_1} (\mathbf{f}_{b_1}, \mathbf{f}_{b_1}), \sum_{b_2} (\mathbf{f}_{b_2}, \mathbf{f}_{b_2}), \dots, \sum_{b_{N_b}} (\mathbf{f}_{b_{N_b}}, \mathbf{f}_{b_{N_b}}) \right],$$

where each block represents the within-biomarker temporal covariance expressed as a negative squared exponential function

$$\Sigma_b(\mathbf{f}_b(t_1), \mathbf{f}_b(t_2)) = \eta_b \exp \left( -\frac{(t_1 - t_2)^2}{2 l_b^2} \right),$$

and where the parameters  $\eta_b$  and  $l_b$  are the marginal variance and length-scale of the biomarker's temporal evolution, respectively.

*Individual random effects.* The random covariance function  $\Sigma_S$  models the individual deviation from the fixed effect, and is represented as a block-diagonal matrix

$$\Sigma_S(\mathbf{v}^j, \mathbf{v}^j) = \text{diag} \left[ \sum_{b_1} (\mathbf{v}_{b_1}^j, \mathbf{v}_{b_1}^j), \sum_{b_2} (\mathbf{v}_{b_2}^j, \mathbf{v}_{b_2}^j), \dots, \sum_{b_{N_b}} (\mathbf{v}_{b_{N_b}}^j, \mathbf{v}_{b_{N_b}}^j) \right],$$

where each block  $\Sigma_b^j$  corresponds to the covariance function associated with the individual process  $\mathbf{v}_b^j(t)$ . Thanks to the flexibility of the proposed generative model, any form of the random effect covariance  $\Sigma_S$  can be easily specified in order to model the subject-specific biomarkers' progression. In what follows we will use a linear covariance form  $\Sigma_b^j(\mathbf{v}_b^j(t_1), \mathbf{v}_b^j(t_2)) = (\sigma_b^j)^2 ((t_1 - \bar{t})(t_2 - \bar{t}))$ , where  $\bar{t}$  is the average observational time for individual  $j$ , when more than 4 measurements are available, and i.i.d. Gaussian covariance form  $\Sigma_b^j(\mathbf{v}_b^j(t_1), \mathbf{v}_b^j(t_2)) = (\sigma_b^j)^2$  when 2 or 3 measurements are available, while assigning it to 0 otherwise (thus by accounting only for the observational noise  $\sigma_b^2$ ). This choice is motivated by stability concerns, in order to keep the model complexity compatible with the generally limited number of measurements available for each individual.

#### Individual time transformation

The generative model (1) is based on the key assumption that the longitudinal observations across different individuals are defined with respect to the same temporal reference. This assumption may be invalid when the temporal alignment of the individual observations with respect to the common group-wise model is unknown, for instance in the typical scenario of a clinical trial in AD where the patients' observational time is relative to the common baseline, and where the disease onset is a latent event (past or future) which is not directly measurable. This modeling aspect is integrated by assuming that each individual measurement is made with respect to an absolute time-frame  $\tau$  through a time-warping function  $t = \phi^j(\tau)$  that models the time-reparameterization with respect to the common group-wise evolution. Model (1) can thus be reparameterized as

$$\mathbf{y}^j(\phi^j(\tau)) = \mathbf{f}(\phi^j(\tau)) + \mathbf{v}^j(\phi^j(\tau)) + \mathbf{e}. \quad (3)$$

The present formulation allows the specification of any kind of time transformation, and in what follows we shall focus on the modeling of a linear reparameterization of the observational time  $\phi^j(\tau) = \tau + \mathbf{d}^j$ . This modeling assumption is mostly motivated by the choice of working with a reasonably limited number of parameters, compatibly with the generally short follow-up time available per individual (cfr. Table 2). Within this setting, the time-shift  $\mathbf{d}^j$  encodes the disease stage associated with the individual relatively to the group-wise model.

#### Monotonic constraint in random-effect multimodal GP regression

Due to the non-parametric nature of Gaussian process regression, we need an additional constraint on model (3) in order to identify a unique solution for the time transformation. By assuming a steady temporal evolution of biomarkers from normal to pathological values, we shall

assume that the biomarker trajectories described by (3) follow a (quasi) monotonic behaviour. This requirement can be implemented by imposing a prior positivity constraint on the derivatives of the GP function. Inspired by (Riihimäki and Vehtari, 2010), we impose a monotonicity constraint by assuming a probit-likelihood for the derivative measurements  $\mathbf{m}(t)$  associated with the derivative process  $\dot{\mathbf{f}}(t) = \frac{d\mathbf{f}(t)}{dt}$  at time  $t$ :

$$p(\mathbf{m}(t) | \dot{\mathbf{f}}(t)) = \Phi \left( \frac{1}{\lambda} \dot{\mathbf{f}}(t) \right), \quad (4)$$

with  $\Phi(z) = \int_{-\infty}^z \mathcal{N}(x|0, 1) dx$ . The quantity  $\lambda > 0$  is an additional model parameter controlling the degree of positivity enforced on the derivative process, with values approaching zero for stronger monotonicity constraint. In what follows, the monotonicity of each biomarker is controlled by placing 10 derivative points equally spaced on the observation domain, and by fixing the  $N_b$  derivative parameters  $\{\lambda_{b_k}\}_{k=1}^{N_b}$  to the value of  $1e-6$ . The position of the derivative points was updated at each iteration, according to the changes of the GP domain.

By following a similar construction, we could equally enforce a monotonic behaviour to the random effects associated with the individual trajectories. This additional constraint would however come with a cumbersome increase of the model complexity, since it would introduce an additional layer of virtual derivative parameters (with associated location) per individual. Moreover, while we are interested in modeling a globally monotonic biomarker trajectory on the fixed parameters, we relax this constraint at the individual level, since some subjects may be characterized by non strictly monotonic time-series due to specific clinical conditions.

*Model likelihood and parameters.* Given the sets of individual biomarker measurements  $\mathbf{y} = \{(\mathbf{y}^j(t_i))_{i=1}^{K_j}\}_{j=1}^N$ , and of  $D$  control derivatives  $\mathbf{m} = \{\mathbf{m}_{b_k}(t')\}_{l=1}^D$  at points  $t' = \{t'_l\}_{l=1}^D$  for the progression of each biomarker  $b_k$ , the random effect GP model posterior is:

$$\begin{aligned} p(\mathbf{f}, \dot{\mathbf{f}}, \mathbf{v}^j | \mathbf{y}, \mathbf{m}) &= \frac{1}{Z} p(\mathbf{f}, \dot{\mathbf{f}} | t, t') p(\mathbf{v}^j | t) p(\mathbf{y} | \mathbf{f}, \mathbf{v}^j) p(\mathbf{m} | \dot{\mathbf{f}}) \\ &= p(\mathbf{f}, \dot{\mathbf{f}} | t, t') p(\mathbf{v}^j | t) p(\mathbf{y} | \mathbf{f}, \mathbf{v}^j) \\ &\quad \prod_k \prod_{l'} \Phi \left( \frac{1}{\lambda} \dot{f}_{b_k}(t'_l) \right), \end{aligned} \quad (5)$$

where  $\mathbf{v} = \{\mathbf{v}^j\}_{j=1}^N$ . Due to the non-Gaussianity of the derivative term  $\Phi$ , the direct inference on the posterior is not possible due to its analytically intractable form. For this reason, we employ an approximate inference scheme based on classical approaches to Gaussian process with binary activation functions (Nickisch and Rasmussen, 2008) (Appendix A.1).

Overall, model (3) is identified by  $(N_j + 3)N_b + N_j$  parameters, represented by the fixed effects and noise  $\theta_G = \{\eta_{b_k}, l_{b_k}, \epsilon_{b_k}\}_{k=1}^{N_b}$ , by the individual random effects parameters  $\theta_G^j = \{\sigma_{b_k}^j\}_{k=1}^{N_b}$  and by the time-shifts  $\mathbf{d}^j$ .

In what follows, the optimal parameters are obtained by maximising the approximated log-marginal likelihood derived from the posterior (5) through conjugate gradient descent, via alternate optimization between the hyper-parameters  $\theta_G$  and  $\theta_G^j$ , and the individuals' time-shifts  $\mathbf{d}^j$ . Regularization is also enforced by introducing Gaussian priors for the parameters  $\theta_G$  and  $\theta_G^j$ .

#### Prediction of observations and individual staging

Gaussian processes naturally allow for probabilistic predictions given the observed data. At any given time point  $t^*$ , the posterior biomarker distribution has the Gaussian form  $p(\mathbf{f}^* | t^*, \mathbf{y}, t, \mathbf{m}, t') \sim \mathcal{N}(\boldsymbol{\mu}^* | \boldsymbol{\mu}^*, \Sigma^*)$ :

$$\boldsymbol{\mu}^* = \Sigma_G(\mathbf{f}(t^*), \mathbf{f}(t)) (\Sigma_{\text{joint}} + \tilde{\Sigma}_{\text{joint}})^{-1} \tilde{\boldsymbol{\mu}}_{\text{joint}} \quad (6)$$

$$\Sigma^* = \Sigma_G(\mathbf{f}(t^*), \mathbf{f}(t^*)) - \Sigma_G(\mathbf{f}(t^*), \mathbf{f}(t)) (\Sigma_{\text{joint}} + \tilde{\Sigma}_{\text{joint}})^{-1} \Sigma_G(\mathbf{f}(t), \mathbf{f}(t^*)), \tag{7}$$

where the matrix  $(\Sigma_{\text{joint}} + \tilde{\Sigma}_{\text{joint}})$  is the joint covariance resulting from the inference scheme detailed in Supplementary Section Appendix A.1 (Riihimäki and Vehtari, 2010).

We also derive a probabilistic model for the individual temporal staging given a set of biomarker observations  $\mathbf{y}^*$ , thanks to the Bayes formula:

$$p(t^*|\mathbf{y}^*, \mathbf{y}, t, \mathbf{m}, t') = p(\mathbf{y}^*|t^*, \mathbf{y}, t, \mathbf{m}, t')p(t^*)/p(\mathbf{y}^*|\mathbf{y}, t, \mathbf{m}, t'), \tag{8}$$

which we compute by assuming an uniform distribution on  $t^*$ , and by noting that  $p(\mathbf{y}^*|t^*, \mathbf{y}, t, \mathbf{m}, t') \sim \mathcal{N}(\boldsymbol{\mu}^*, \Sigma^* + \Sigma_\epsilon)$ . In particular, the covariance form  $\Sigma_G(\mathbf{f}(t^*), \mathbf{f}(t^*))$  can be specified in order to account for incomplete data, and thus generalizes the GP model for predictions in presence of *missing biomarker observations*. The posterior distribution (8) quantifies the *confidence* of the model about the individual disease staging, and thus is a valuable information about the precision of the diagnosis. We will also compute the *expectation* of the distribution  $p(t^*|\mathbf{y}^*, \mathbf{y}, t, \mathbf{m}, t')$ , which provides a scalar value that can be used in subsequent classification methods.

**Materials and methods**

**Study participants**

Data used in the preparation of this article were obtained from the ADNI database (<http://adni.loni.usc.edu>). The ADNI was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W. Weiner, MD. The primary goal of ADNI has been to test whether serial magnetic resonance imaging, positron emission tomography, other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of MCI and early AD. For up-to-date information, see [www.adni-info.org](http://www.adni-info.org).

**Data processing**

We collected longitudinal measurements for the ADNI individuals with baseline values of cerebrospinal fluid (CSF)  $A\beta$  amyloid lower than the nominal value of 192 pg/ml. The information was extracted from the ADNIMERGE<sup>2</sup> R package (R Core Team, 2015) (MEDIAN field of the UPENNBIOBK\_MASTER table). This preliminary selection is aimed to validate the model on a clinical population likely to represent the whole disease time-span.

The model was trained on a group of 200 randomly selected individuals including healthy volunteers, mild cognitive impairment subjects converted to AD (MCI conv), and AD patients having at least one measurement for each of the following biomarkers: *volumetric measures* (hippocampal, ventricular, entorhinal, and whole brain volumes), *glucose metabolism* (average normalized FDG uptake in prefrontal cortex, anterior cingulate, precuneus and parietal cortex), *brain amyloidosis* (average normalized AV45 uptake in frontal cortex, anterior cingulate, precuneus and parietal cortex), and *functional, neuropsychological and cognitive function* measured by common scores (ADAS13, RAVLT learning, and FAQ).<sup>3</sup> The testing set was composed of the remaining 582 subjects, including a subgroup of MCI non converted to AD during the observational time (MCI stable). The image-derived measures used in the study (volumetric MRI and average uptake values for AV45- and FDG-PET) were the scalar estimates reported in the ADNIMERGE package (adni-merge table). The volumetric measures were scaled by the individual

total intracranial volume, and all the biomarkers measurements were converted into quantile scores (0–1 for normal to abnormal values), with respect to the biomarkers distribution of the *training set*. This latter modeling precaution is aimed to avoid spurious correlation between training and testing data due to the combined normalization of the values.

The modeling results were evaluated with respect to the baseline diagnostic information reported in the ADNI database, assessed according to the WMS and NINCDS/ADIRDA AD criteria (McKhann et al., 1984). Conversion to MCI or AD was established according to the last follow-up information. Moreover, the MCI group was composed by 138 individuals with baseline diagnosis of early MCI, assessed through the Wechsler Memory Scale Logical Memory II. Among these subjects, 14 of them were in the training group (26% of the total MCI training set size), while the remaining 124 were in the testing set (35% of the total MCI testing set size).

Table 1 shows baseline clinical and sociodemographic information of the individuals used respectively in training and testing set, while in Table 2 we report the average follow-up time and the ratio of missing data of the pooled sample. Supplementary Section Appendix A.2.6

**Table 1**  
Baseline sociodemographic and clinical information for training and testing study cohort. NL: normal individuals, MCI: mild cognitive impairment, AD: Alzheimer's patients. ADAS13: Alzheimer's Disease Assessment Scale-cognitive subscale, 13 items; FAQ: Functional Assessment Questionnaire; RAVLT learning: Rey Auditory Verbal Learning Test, learning item; FDG: (18)F-fluorodeoxyglucose positron emission tomography (PET) imaging; AV45: (18)F-florbetapir Amyvid PET imaging.

Group	NL	NL converted	MCI stable	MCI converted	AD
<b>Training data</b>					
N	67	5	0	53	75
Age	73 (6)	81.4 (5.2)	/	72 (7.7)	73 (8.5)
Sex (% females)	61	0	/	43	45
Education (yrs)	16.2 (2)	17.2 (3)	/	15.8 (2.6)	16 (2.4)
ADAS13	8.8 (4.5)	13.8 (2.4)	/	22.6 (6.7)	31.3 (8.5)
FAQ	0.2 (0.6)	0.4 (0.5)	/	5.2 (4.5)	12.9 (7)
RAVLT learning	5.6 (2.6)	2.2 (1.9)	/	3.2 (2.5)	1.8 (1.7)
Entorhinal (cm <sup>3</sup> )	3.9 (0.6)	3.7 (0.5)	/	3.2 (0.7)	2.9 (0.6)
Hippocampus (cm <sup>3</sup> )	7.5 (0.9)	6.7 (0.7)	/	6.2 (0.9)	6 (9.3)
Ventricles (cm <sup>3</sup> )	36 (20)	57 (26)	/	42 (21)	47 (22)
Whole brain (cm <sup>3</sup> )	1057 (105)	1106 (116)	/	1040 (107)	1013 (113)
FDG	6.6 (0.5)	6.1 (0.65)	/	5.7 (0.6)	5.2 (0.64)
AV45	1.2 (0.2)	1.3 (0.09)	/	1.4 (0.2)	1.4 (0.2)
<b>Testing data</b>					
N	74	17	243	106	145
Age	75.3 (5.9)	76.5 (4)	73.3 (7)	73.6 (7.3)	75 (7.9)
Sex (% females)	55	41	39	40	39
Education (yrs)	16 (2.9)	16.2 (2.6)	16 (2.8)	16 (3)	15.3 (3.1)
ADAS13	9.8 (4)	11.7 (3.4)	15.7 (6)	21 (6.1)	29.4 (8.2)
FAQ	0.5 (1.3)	0.6 (1.6)	2.7 (3.5)	5.1 (4.7)	12.9 (6.8)
RAVLT learning	5.6 (2.2)	5.6 (2.7)	4.3 (2.5)	2.8 (2.2)	1.8 (1.9)
Entorhinal (cm <sup>3</sup> )	3.8 (0.4)	3.6 (0.7)	3.6 (0.7)	3.1 (0.7)	2.7 (0.7)
Hippocampus (cm <sup>3</sup> )	7.2 (0.7)	7.2 (0.8)	6.9 (1)	6 (0.8)	5.7 (0.1)
Ventricles (cm <sup>3</sup> )	33 (15)	44 (21)	39 (23)	41 (23)	49 (24)
Whole brain (cm <sup>3</sup> )	1019 (102)	1055 (93)	1056 (100)	992 (110)	972 (124)
FDG	6.5 (0.62)	6.4 (0.7)	6.3 (0.7)	5.9 (0.6)	5.4 (0.7)
AV45	1.21 (0.19)	1.4 (0.2)	1.3 (0.19)	1.4 (0.2)	1.4 (0.2)

<sup>2</sup> [adni.bitbucket.io/adnimerge.html](http://adni.bitbucket.io/adnimerge.html).

<sup>3</sup> ADAS13: Alzheimer's Disease Assessment Scale-cognitive subscale, 13 items; FAQ: Functional Assessment Questionnaire; RAVLT learning: Rey Auditory Verbal Learning Test, learning item; FDG: (18)F-fluorodeoxyglucose positron emission tomography (PET) imaging; AV45: (18)F-florbetapir Amyvid PET imaging.

**Table 2**  
Average follow-up years and percentage of individuals with missing data (in parenthesis).

Ventr	Hippo	Ent	Whole Brain	ADAS13	FAQ	RAVLT	AV45	FDG
Training data								
2.3 (0)	2.3 (0)	2.3 (0)	2.3 (0)	3 (0)	3.3 (0)	3.3 (0)	1.9 (0)	1.6 (0)
Testing data								
3.4 (11)	3.4 (11)	3.4 (11)	3.4 (11)	3.9 (0)	3.9 (0)	3.9 (0)	3.8 (43)	3 (19)

reports the R code used for data pre-processing.

### Longitudinal modeling of Alzheimer's disease progression

#### Model training

The model was applied in order to estimate the temporal biomarker evolution and the disease stage associated with each individual in the training set. The plausibility of the model was assessed by group-wise comparison of the predicted time-shift, and by correlation with respect to the time to AD diagnosis for the MCI individuals subsequently converted to AD. For sake of comparison we also correlated the progression modelled with our approach with respect to the one estimated with the method proposed in (Donohue et al., 2014). The method was applied to the training data by using the standard parameters defined in the R package *GRACE*<sup>4</sup> (see Supplementary Material Appendix A.2.2 for further details).

#### Model testing on de-novo individuals

The estimated probabilistic disease progression model provides a valuable clinical reference, as it can be used to predict an individual pathological stage, as well as to quantify the biomarkers predictive value, or the influence of missing data. To this end, we estimated the predictive performance of the model in assessing the individual pathological stage with respect to follow-up assessments and missing biomarkers. This was done by estimating the predictive accuracy of the group-wise separation obtained via increasing thresholds of the estimated temporal progression.

## Results

### Model plausibility

The estimated biomarker progression (Fig. 1-A) shows a biologically plausible description of the pathological evolution, compatible with previous findings in longitudinal studies in familial AD (Bateman et al., 2012), and with the hypothetical models of AD progression (Jack et al., 2010; Frisoni et al., 2010). The progression is defined on a time scale spanning roughly 20 years, and is characterized at the initial stages by high-levels of AV45, followed by the abnormality of ventricles volume, of FDG uptake, and of the whole brain volume. These latter measures are however heterogeneously distributed across clinical groups, and with rather large variability. The evolution is further characterized by increasing abnormality of the volumetric measures (especially hippocampal volume), and by the steady worsening of neuropsychological scores such as FAQ. The model thus shows that the transition from normal to pathological levels is essentially characterized by increase of hypo-metabolism, followed by the pronounced temporal brain atrophy. Moreover, the worsening of the neuropsychological and functional scores closely (almost linearly) follows the progression in the advanced clinical stages. The joint visualization of the temporal progression of the biomarkers with temporal derivative of the modelled average trajectories is shown in Supplementary figure A10. The illustration confirms that ADAS13 and FAQ are characterized by very similar longitudinal profiles, and show the largest changes in the latest stages of the pathology (peak of the derivative at  $t > 0$ ). On the contrary, the change in hippocampal

volume is more strongly associated with the earlier stages of the pathology. AV45 and ventricles volumes are the least informative and are associated with the lowest changes.

Fig. 1-B (top) shows the posterior time-shift distributions associated with the individuals. The distributions denote the confidence of the model in associating to each individual a temporal staging with respect to the global pathological progression. The boxplot of Fig. 1-B (bottom) reports the group-wise expectation of the individual time-shifts. Healthy individuals (blue) are associated with the early stages of the pathology in both training and testing data, while MCI (purple) and AD patients (red) are characterized by respectively intermediate and late predicted progression stages. The group-wise comparison between the expected time-shifts was statistically significant between each group pairs (ANOVA,  $p < 1e-6$ ). Moreover, the time to conversion to AD in the MCI group was significantly correlated with the disease staging quantified by the expectation of the individual time distributions ( $R^2 = -0.4$ ,  $p = 3.8e - 4$ ).

Finally, when applying (Donohue et al., 2014) to the training data we measured a strong agreement between the resulting progression and the one obtained with our method, resulting in a correlation between the corresponding individual time-shifts of 0.94 ( $p < 1e-6$ ) (Supplementary Material Appendix A.2.2).

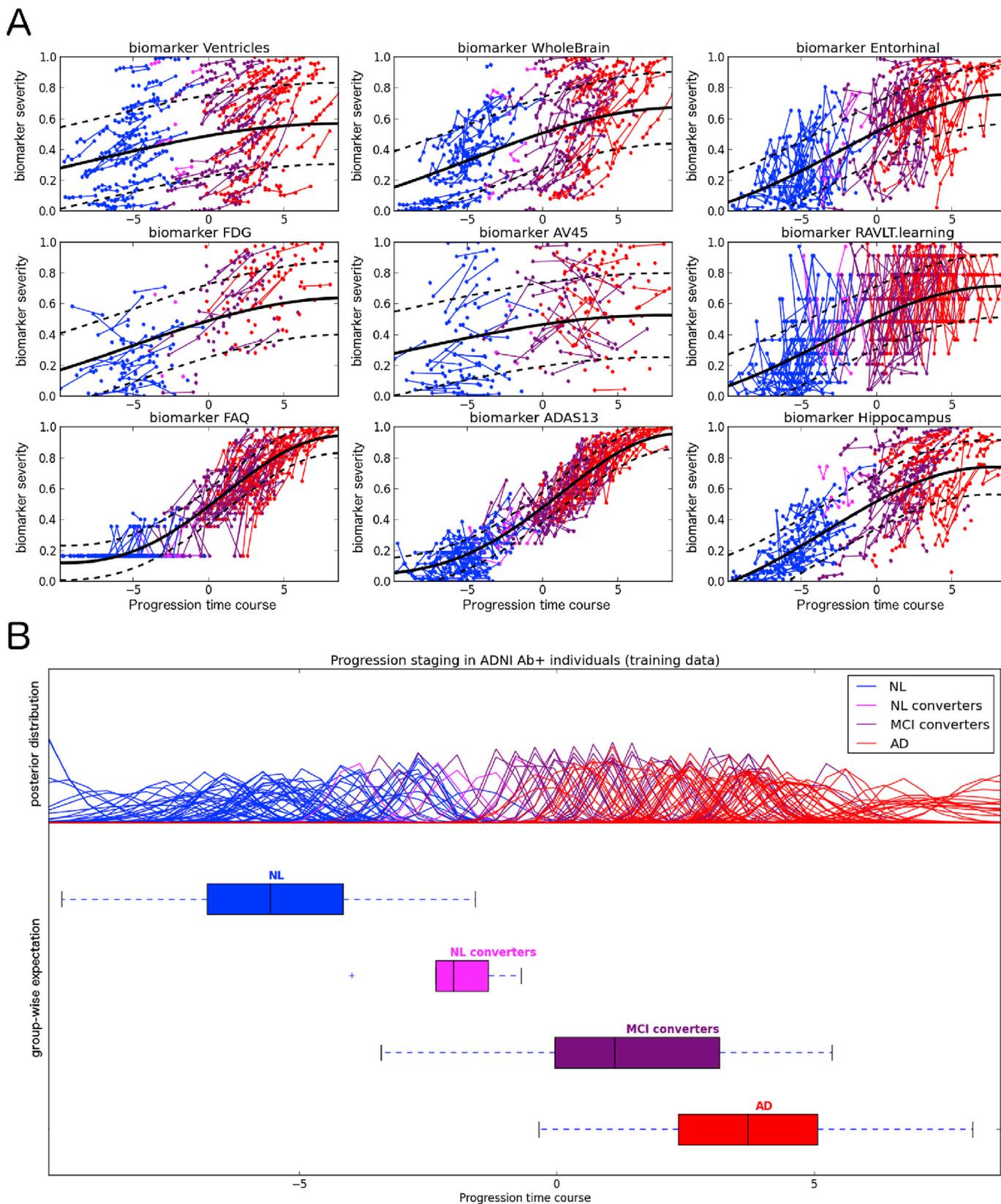
### Assessing diagnostic uncertainty in testing data: an illustrative example

This section illustrates the use of the model represented in Fig. 1 for the quantification of diagnostic uncertainty in testing individuals. We consider the hypothetical scenario where the baseline values for different biomarkers are measured for a given patient, namely FAQ, hippocampal and ventricle volumes. We assume that the biomarkers values correspond to the 20th percentiles with respect to the biomarkers distribution of the training set (i.e. FAQ = 1, normalized Hippo =  $5e-3$ , normalized Ventr =  $1.7e-2$ ). Fig. 2 (left) shows the disease staging prediction obtained with formula (8) based on the value of each biomarker. We note that FAQ and hippocampal volume lead to similar posterior Gaussian distributions of disease staging, with expectation of respectively  $t_{FAQ} = -6$  and  $t_{hippo} = -5.6$  (indicated by the vertical lines in the figure), and standard deviation of  $sd_{FAQ} = 6.3$  and  $sd_{hippo} = 5.9$ . The prediction associated with ventricles volume is wider and associated with higher uncertainty, with mean and standard deviation of respectively  $t_{ventr} = -3.8$  and  $sd_{ventr} = 6.1$ .

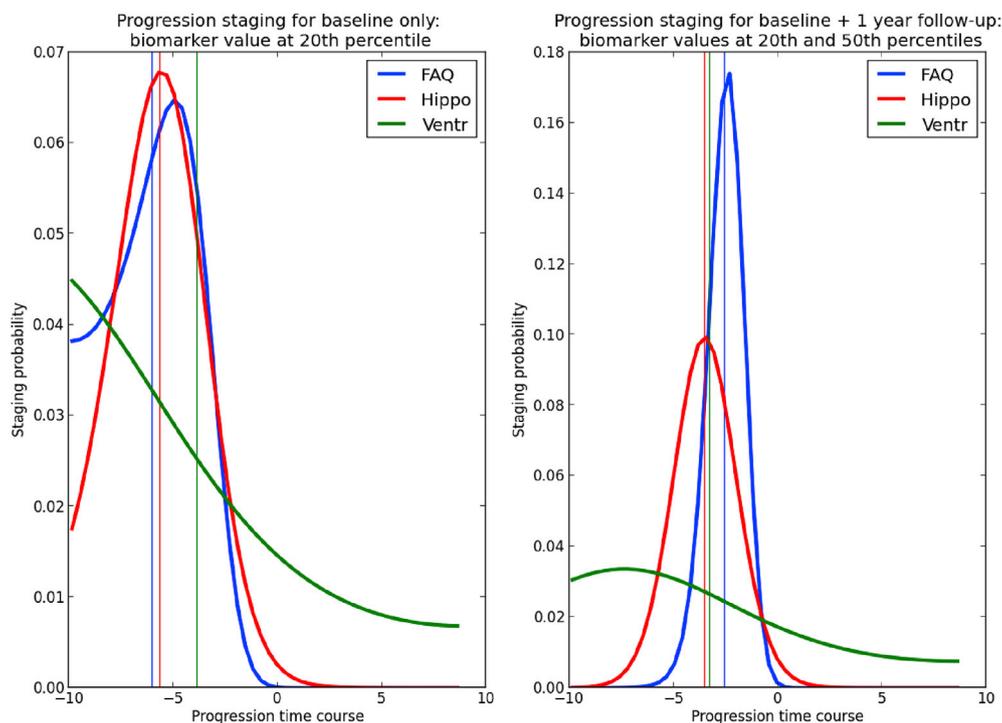
We now suppose that for the same patient we acquire a follow-up measurement for each biomarker at year 1, with values corresponding to the 50th percentiles of the distribution of the training set (i.e. FAQ = 5, normalized Hippo =  $4.3e-3$ , normalized Ventr =  $2.7e-2$ ). The right hand side of Fig. 2 shows the new prediction based on the joint baseline + follow-up information. For each biomarker the posterior distributions indicate an increase of the predicted disease stage with respect to the baseline scenario, while the prediction uncertainty is generally lower. Although the expectation for the 3 biomarkers is very similar ( $t_{FAQ} = -2.5$ ,  $t_{hippo} = -3.5$ , and  $t_{ventr} = -3.2$ ), we notice that FAQ leads to the highest diagnostic confidence ( $sd_{FAQ} = 2.6$ ), followed by hippocampal volume ( $sd_{hippo} = 3.8$ ), and finally by ventricles volume ( $sd_{ventr} = 5.7$ ). Further assessment of the relationship between biomarker variability and model prediction is provided in supplementary Section Appendix A.2.3.

This illustrative example shows that the proposed probabilistic

<sup>4</sup> <https://mdonohue.bitbucket.io/grace/>.



**Fig. 1.** A) Modelled biomarker progression in the training set of 200  $A\beta$  amyloid positive individuals (solid/dashed lines: mean  $\pm$  sd). B) Posterior prediction for the individual time shift in training data (top: individual time-shift distribution; bottom: group-wise boxplot of the expected time-shift). Healthy individuals are generally displaced at the early stages of the pathology, while the predictions for MCI and AD patients are associated with respectively intermediate and late progression stages. NL: normal individuals, MCI: mild cognitive impairment, AD: Alzheimer's patients.



**Fig. 2.** Illustrative example: posterior prediction of disease staging for a testing individual based on baseline (left), and baseline + follow-up (right) information for three biomarkers: FAQ, hippocampal and ventricles volume. The biomarker values correspond to the 20th and 50th percentiles of the training-group distribution for respectively baseline and follow-up measures. Adding the follow-up information leads to increased estimates of the disease staging and to generally lower prediction uncertainty. Although the distributions associated with different biomarkers generally lead to similar expectations, FAQ and hippocampal volume lead to the lowest diagnostic uncertainty. Vertical lines: expectation for each posterior distribution.

framework represents a valuable instrument for the assessment of the diagnostic value and uncertainty associated with different biomarkers, and can faithfully track the pathological progression of testing individuals along the modelled trajectories, from normal to pathological levels.

#### DPM for probabilistic diagnosis in ADNI

We now assess the predictive results of the model when applied to the testing ADNI cohort. Fig. 3 shows the individual posterior predictive distributions associated with the testing individuals, and the boxplot of the expected time-shift when using the model as statistical reference through formula (8). The figure reports the two different modeling scenarios based on baseline information only (Fig. 3–1), and on the complete set of baseline and longitudinal measurements (Fig. 3–2). We first note that the group-wise differences between the expected time-shifts are compatible for both scenarios, as shown by the similar boxplot distributions across groups (Fig. 3-1b vs 3-2b). The consistency of the predictions is further illustrated in Figure A.9, where it is shown that the group-wise distribution and ordering of the predicted time-shifts in the testing data are compatible with those estimated in the training one.

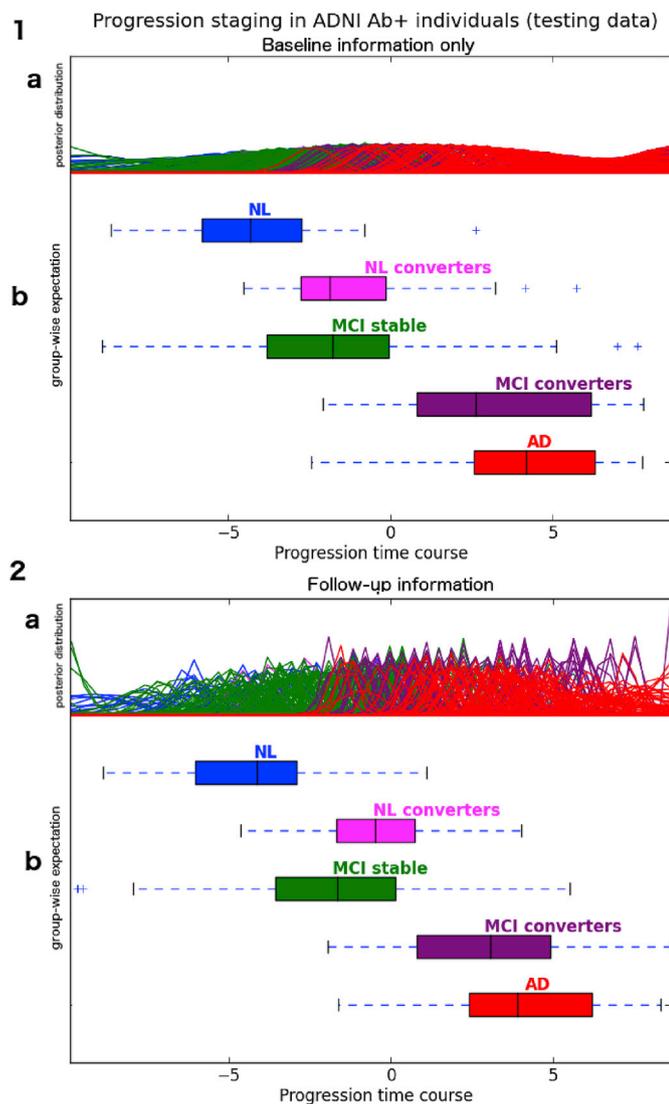
However, the joint use of baseline and follow-up information largely reduces the uncertainty of the predictions (Fig. 3-1a vs 3-2a). Indeed, the time distributions predicted when using baseline and follow-up information are narrower as compared to the wider confidence margins obtained by using the baseline information only. Therefore, adding follow-up measurements importantly improves the confidence of the model in determining the individual pathological stage.

As with the training case, for both scenarios the group-wise distribution of the expected time-shift shows a significant separation between the clinical groups according to the increase of the pathological stage (ANOVA,  $p < 1e-4$ ). Interestingly, the temporal positioning of the non converting MCI lies between controls and MCI converters, and is on average lower than the one of healthy individuals subsequently converted to cognitive impairment.

Fig. 4 reports the classification results based on the baseline information only, and on increasing thresholds of the progression time course. Although the model is not optimized to explicitly classify the clinical groups, the simple thresholding based on the model predictions generally shows high face validity with respect to the clinical diagnosis. For all the considered scenarios, the highest accuracy is reached in a time window around the point  $t = 0$ , while the area under the ROC curve is 0.99, 0.88 and 0.83 for NL vs AD, MCI converters vs MCI stable, and NL converters vs NL stable, respectively.

We further tested the model in presence of missing information, by computing the predictions when only one baseline biomarker is available (Fig. 5). The predictive outcomes show important variations depending on the considered biomarker, while the confidence bounds for the predictions are usually large, to denote increased uncertainty. We also note that FAQ, ADAS13, and hippocampal volume are the biomarkers leading to the largest group-wise separation, along with the lowest prediction uncertainty. This aspect is quantified in Table 3, reporting the discrimination results with respect to the nominal cut-off point of  $t = 1.65$ , corresponding to the 15th percentile of the distribution of the expected time-shift in the training AD group, as well as the area under the receiving operating characteristic curve (AUC). Although the highest discriminative results are consistently obtained when the biomarkers are used jointly, we note that the neuropsychological tests generally lead to the best predictive performance in identifying AD patients with respect to healthy individuals, followed by brain hypo-metabolism (FDG-PET), and temporal atrophy (Entorhinal and Hippocampal volume). This is related to the lower uncertainty of the modelled progressions, which leads to a more accurate identification of the individual staging along the pathological trajectory. The scenario sensibly changes in the other comparison scenarios (MCI conv vs stable and NL conv vs stable), where the sensitivity of the neuropsychological scores shows an important drop, while the other biomarkers (especially hippocampal and entorhinal volumes) provide comparable or even better discriminative performances.

These figures were similar when considering the single biomarkers within the longitudinal setting, where the neuropsychological tests still



**Fig. 3.** Posterior prediction for the individual time shift in testing data by using i) only the baseline information (1a-b), and ii) the baseline + follow-up information available for each testing subject (2a-b). Healthy individuals are generally displaced at the early stages of the pathology, while the predictions for MCI and AD patients are associated with respectively intermediate and late progression stages. The results are similar for both scenarios, although by adding the follow-up information we largely reduce the uncertainty in the prediction of the individual's pathological stage (subfigure 1a vs 2a). NL: normal individuals, MCI: mild cognitive impairment, AD: Alzheimer's patients.

outperformed the other biomarkers in discriminating the clinical groups (Supplementary Figure A11).

For the sake of comparison we finally benchmarked the predictive results provided by the disease progression model with respect those obtained by the classification analysis performed with standard statistical tools, such as a random forest classifier. We note that the comparison of the classification performance obtained on the heterogeneous data considered in this work is generally challenging, since the proposed DPM 1) accounts for missing data and non-fixed number of time points per individuals, and 2) is formulated in order to consistently handle both longitudinal and cross sectional measurements, either for training and prediction. To date there is not a consensus on the optimal approach to adopt to tackle these important modeling constraints, while the comparison between the classification performance obtained with complex machine learning methods is currently matter of scientific debate and investigation (Mendelson et al., 2016).

For this reason we restricted the random forest classification task to a standard statistical setting, in order to essentially provide a reliable benchmark for the classification performance of the proposed disease progression. To this end we trained the random forest on the classification between healthy individuals and AD patients based on the baseline measurements of the training group, while the missing entries in the testing data were imputed via nearest neighbour search, based on the available biomarkers. The classification results are reported in Supplementary Table A5.

The performance of the random forest classifier is generally inferior to the one obtained with the proposed approach, as witnessed by the consistently lower AUC obtained for all the comparisons. The difference becomes more evident for the more challenging classification problems, such as the identification of conversion in MCI and healthy individuals. This result is indicative of the reliability of the classification results obtained by the proposed disease progression model, especially when considering that the random forest classifier is explicitly optimized to maximize the separation between groups, while the accuracy results reported in Table 3 are based on the empirical choice of a reference threshold in the training population.

#### DPM staging and chronological age

We finally compare the relationship between the predicted disease staging in training and testing set and the individual chronological age. We first note that both training and testing clinical groups were matched by age, with the exception of the 5 training healthy subjects converted to MCI (or AD) that were slightly older with respect to the reference training healthy population ( $p = 0.02$ ).

Nevertheless, when comparing the estimated time shift with respect to the chronological age of each individual we didn't report any significant correlation between these measures. Interestingly, the same lack of association is also quantifiable in the testing group (Fig. 6). This result, in association with the strong relationship between time shift and clinical condition reported in Section 3.3, let us conclude that the model is describing the biomarker's variation essentially related to the pathological progression, which is orthogonal to the effect of healthy aging quantified by the chronological age. This result points to the effectiveness of the proposed approach in capturing significant effects related to the specific temporal progression of the disease.

#### Discussion

This study explores the use of DPM for probabilistic diagnosis and uncertainty quantification in an hypothetical clinical scenario. The proposed approach is based on the reformulation of DPM through a novel probabilistic approach aimed at leveraging on the longitudinal modeling of disease progression for prediction and quantification of the diagnostic uncertainty in neurodegeneration, by optimally combining the information provided by the several biomarkers into a biologically plausible and intelligible score quantified by the time shift. This work thus extends the previous contributions by proposing DPM as a probabilistic tool for diagnostic purposes, which can be used to quantify staging and predictive uncertainty of de-novo individuals in clinical trials. The disease progression model itself thus can be seen as a novel biomarker of pathological progression. We also note that the time shift is a *relative* measure of disease progression accounting for the biomarker variability observed in the training population. Thus, the point 0 is generally not associated with the conversion to AD, as it is relative to the data initialization (in this case the study baseline).

We illustrated the use of DPM as benchmarking tool for the statistical comparison of biomarkers. The model allows the quantification of the variability associated with the single biomarkers, by identifying the related uncertainty in characterizing the progression from normal to pathological levels. The proposed model can be thus used as a reference for screening and enrichment purposes in clinical trials (Lorenzi et al.,

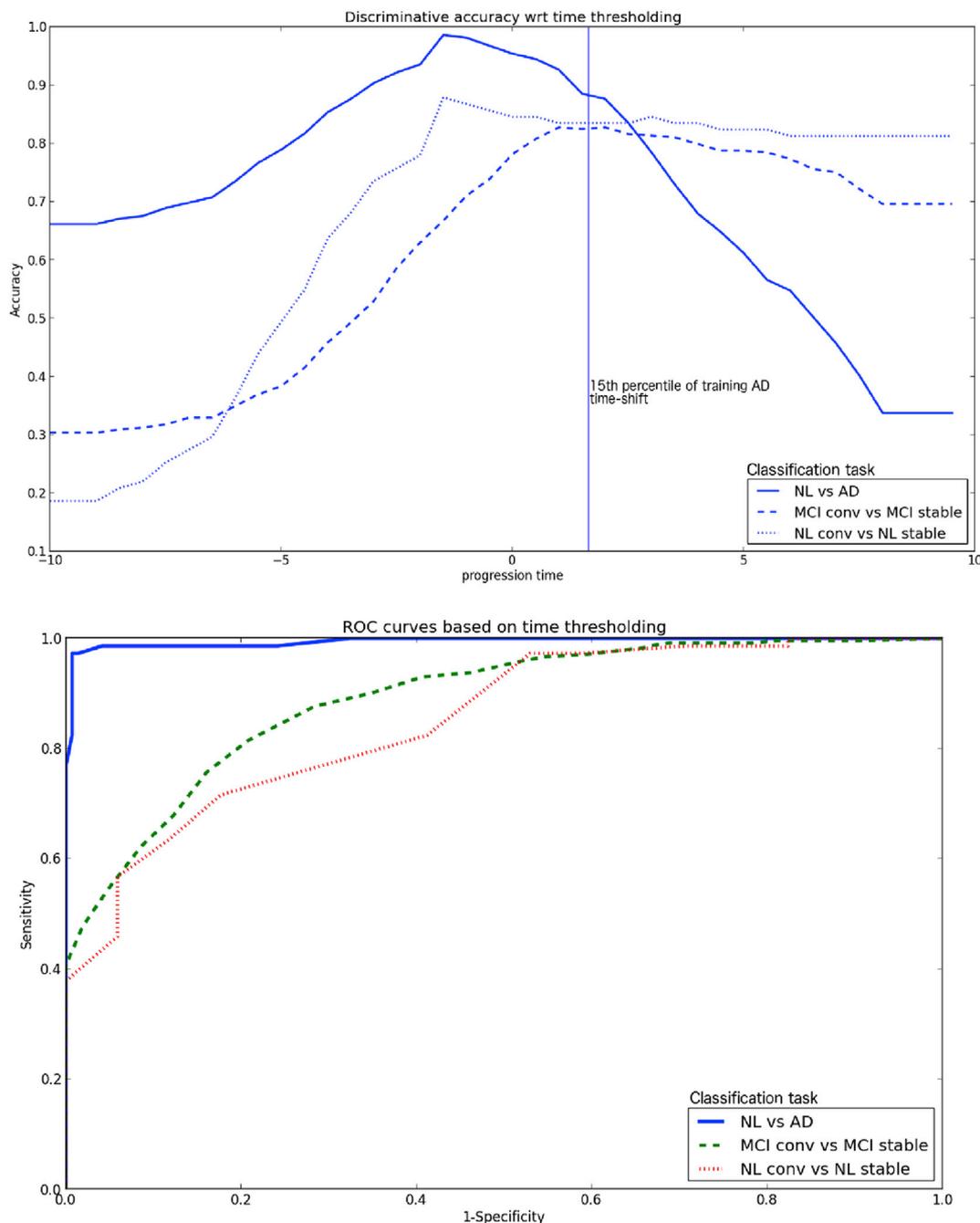


Fig. 4. Predictive accuracy of the model when considering the joint set of available biomarkers measurements. The vertical bar indicates the reference threshold value of  $t = 1.65$ , corresponding to the 15th percentile of the time distribution of the training AD group. MCI: individuals with mild cognitive impairment, AD: Alzheimer's patients.

2010; Yu et al., 2014; Hill et al., 2014).

The modelled progression showed that neuropsychological tests generally lead to lower uncertainty for identifying the individual clinical stage, and to the higher separation power between healthy and AD groups. This finding is compatible with the results reported by previous disease progression models applied to ADNI, such as (Jedynak et al., 2012) and (Young et al., 2014). In this latter study ADAS13 consistently appeared among the first events distinguishing the normal disease stages from the pathological ones. Furthermore, our analysis further showed that volumetric measures such as hippocampal and entorhinal volumes provide equivalent if not superior diagnostic performances when tested on the more challenging problem of detecting conversion to dementia from healthy and MCI stages, especially in terms of improved AUC. Nevertheless, some care should be taken in drawing conclusions from the

present analysis. Our model was based on the standard volumetric measures provided in the ADNI database, and we cannot exclude that a more precise quantification of morphological brain changes would lead to even better performance of volumetric biomarkers (Wolz et al., 2010; Cash et al., 2015). Furthermore, the proposed model was not optimized in order to maximize the classification accuracy between clinical groups. For example, the results reported in Table 3 are based on the choice of the temporal threshold corresponding to the reference value of the 15th percentile of the AD distribution. This cut-off was not optimized to maximize the predictive outcome of the biomarkers, but was rather chosen based on heuristics aimed at illustrating the use of the model for predictive purposes. We thus cannot exclude that the optimization of the temporal threshold would lead to different figures for the classification task. The reported results are therefore indicative of the effectiveness of

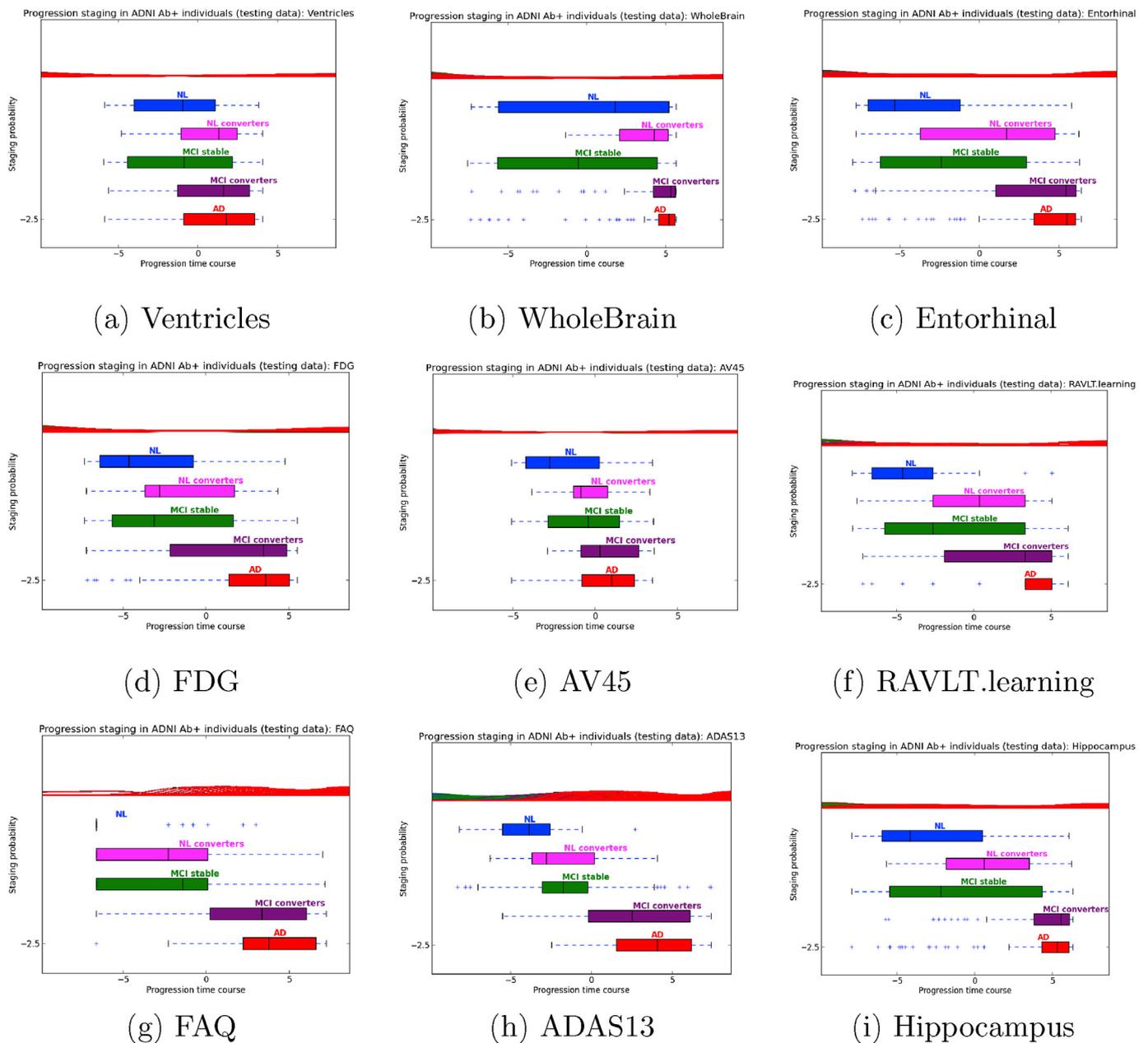


Fig. 5. Posterior prediction on testing data by using a single biomarker and the baseline information only.

Table 3

Classification results by using the reference time threshold of  $t = 1.65$ , corresponding to the 15th percentile of the training AD time distribution.

	Biomarker									
	all	Hippo	Ventr	WholeBr	Entor	FDG	AV45	RAVLT	FAQ	ADAS13
NL vs AD (145 vs 74)										
Accuracy	0.89	0.81	0.62	0.76	0.83	0.80	0.63	0.82	0.88	0.83
Sensitivity	0.83	0.84	0.52	0.9	0.82	0.74	0.82	0.76	0.84	0.75
Specificity	0.98	0.76	0.80	0.46	0.83	0.89	0.46	0.94	0.97	0.98
AUC	0.99	0.87	0.69	0.7	0.89	0.87	0.73	0.91	0.98	0.98
MCI conv vs MCI stable (106 vs 243)										
Accuracy	0.82	0.67	0.62	0.69	0.7	0.71	0.69	0.67	0.79	0.79
Sensitivity	0.65	0.85	0.5	0.89	0.74	0.65	0.37	0.56	0.63	0.54
Specificity	0.90	0.59	0.68	0.60	0.68	0.73	0.75	0.71	0.86	0.9
AUC	0.88	0.79	0.61	0.78	0.76	0.74	0.61	0.66	0.81	0.82
NL conv vs NL stable (17 vs 74)										
Accuracy	0.83	0.70	0.71	0.54	0.77	0.76	0.73	0.83	0.82	0.83
Sensitivity	0.18	0.47	0.41	0.82	0.52	0.29	0.27	0.35	0.17	0.17
Specificity	0.98	0.77	0.80	0.47	0.83	0.89	0.86	0.94	0.97	0.98
AUC	0.83	0.71	0.65	0.63	0.74	0.65	0.65	0.7	0.63	0.68

## Chronological age vs GP model staging

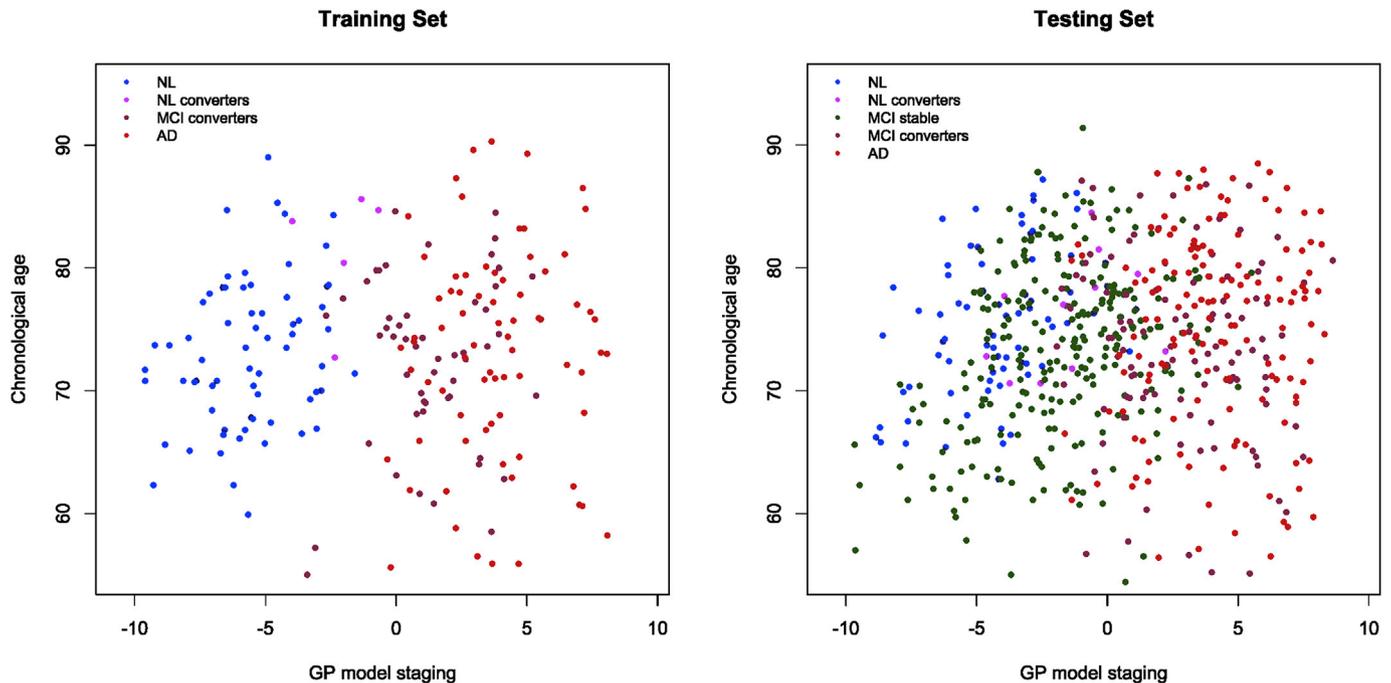


Fig. 6. Chronological age (y-axis) vs model staging (x-axis). The estimated time-shift is decorrelated from the chronological age in both training and testing data ( $p > 0.05$ ).

the model in faithfully representing the clinical spectrum of the disease. We note also that the reported figures are in line with those provided by state-of-art methods in AD classification, without requiring complex parameter optimization procedures, which would introduce additional levels of cross-validation and expose the results to selection bias and generalization issues (Mendelson et al., 2016).

Thanks to the probabilistic formulation we showed that the use of longitudinal information is important for reducing the uncertainty of the prediction, and thus allowing one to better identify the disease status associated to an individual. This important aspect is in agreement with the generally higher statistical power reported in previous Alzheimer's studies comparing longitudinal measurements to baselines ones (Heneman et al., 2009; Frisoni et al., 2010; Xu et al., 2014).

In this work we focused on the modeling of the progression of amyloid positive individuals. This choice was motivated by the interest in assessing the model performance on an homogeneous clinical population likely to be representative of the Alzheimer's evolution. While the absence of pathological amyloid levels seems indicative of non-Alzheimer's pathophysiology (Gordon et al., 2016; Mormino et al., 2016), there is currently an active debate on the mechanisms of neurodegeneration not related to brain amyloidosis (Jack et al., 2016). The investigation of these aspects goes beyond the scope of the present work, and future extensions of disease progression modeling will aim at identifying differential progressions underlying sub-pathologies, for example by reformulating the proposed random effect regression within the realm of Gaussian process mixture models (Lázaro-Gredilla et al., 2012; Ross and Dy, 2013). Analogously, the MCI population used for model training was composed exclusively by MCI individuals subsequently converted to AD, in order to train the model on a homogeneous data most likely to include the largest representation of individuals effectively affected by Alzheimer's disease. Although the inclusion of MCI stable could provide additional information on the intermediate pathological stages, this choice may probably lead to larger variability in the training set, as stable MCI are generally characterized by larger heterogeneity, either cross-sectionally and longitudinally, and higher diagnostic uncertainty. This modeling choice was also motivated by practical reasons since,

thanks to the adopted data selection scheme, we were able to validate the model on a *large and independent* set of testing individuals including an important sample of MCI individuals across different clinical stages, thus providing a thorough and stringent assessment of the predictive qualities of the proposed approach.

#### Methodological considerations

From the methodological perspective, we proposed a novel probabilistic approach based on Gaussian process regression for disease progression modeling from time-series of biomarker measurements enabling novel applications beyond the state-of-art, such as the probabilistic prediction of disease staging in testing individuals. Furthermore, the model naturally accounts for missing data, and provides uncertainty quantification of the biomarker evolutions. Similarly to (Donohue et al., 2014), in this work we focused on the modeling of disease staging represented by a time shift, although the proposed framework can naturally account for more complex time transformations, provided that a sufficient number of time points is available for each individual.

From the methodological point of view, the proposed model extends current approaches to GP-regression by consistently integrating time-reparameterization and monotonic constraints within a random effect regression framework. Monotonic GPs were introduced in (Riihimäki and Vehtari, 2010) as a principled regularization solution to improve the plausibility of modeling results. For example, the strength of such a regularization approach in biomedical applications has been illustrated in survival analysis (Joensuu et al., 2012). Our approach extends this framework by consistently integrating a latent time variable parameter within a random effect model formulation.

The idea of estimating a time transformation in a GP regression framework has been previously used by (Liu et al., 2010) to account for uncertain measurement times to a microarray dataset of mRNA. However, in that work the estimation of the time uncertainty was subject to a strong prior constraint based on the assumption that the unknown biological time must be similar to the measured one. In the application proposed in our work such an assumption is no longer valid and would

ultimately lead to implausible estimations. On the contrary, the proposed GP regression is able to recover the underlying time transformation thanks to the proposed monotonicity regularization.

Finally, thanks to the flexibility of the proposed Gaussian process framework, further extensions of the model will enable to consistently integrate a spatio-temporal covariance model, such as the efficient Kronecker form of (Lorenzi et al., 2015), to provide a unified framework for jointly modeling time series of images and scalar biomarkers data in a coherent fully Bayesian setting.

## Conclusions

This work proposes an extension of DPM for the accurate quantification of the diagnostic uncertainty in Alzheimer's disease. The proposed application shows that DPM provides at the same time a plausible description of the transition from normal to pathological stages along the natural history of the disease, as well as remarkable diagnostic performances when tested on de-novo individuals. The model used in this study can account for any missing data patterns (longitudinal or across biomarkers), and allows to directly quantify the uncertainty related to the missing information. It thus represents a novel and promising tool for the analysis of clinical trials data.

## Further information

The open-source code as well as the proposed predictive model trained on ADNI data will be available at the author's web-page: <https://team.inria.fr/asclepios/marco-lorenzi/>. The realization of this study required about 1.5kWh of computing power.

## Acknowledgments

EPSRC grants EP/J020990/01 and EP/M020533/1 support DCA and SO's work on this topic. DCA and SO also received support from the European Union's Horizon 2020 research and innovation programme under grant agreement No 666992 (EuroPOND) for this work. MF gratefully acknowledges support from the AXA Research Fund.

Data collection and sharing for this project was funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer's Association; Alzheimer's Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Eisai; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health ([www.fnih.org](http://www.fnih.org)). The grantee organisation is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer's Disease Cooperative Study at the University of California, San Diego. ADNI data are disseminated by the Laboratory for NeuroImaging at the University of Southern California.

## Appendix A. Supplementary data

Supplementary data related to this article can be found at <https://doi.org/10.1016/j.neuroimage.2017.08.059>.

[org/10.1016/j.neuroimage.2017.08.059](https://doi.org/10.1016/j.neuroimage.2017.08.059).

## References

- Bateman, R.J., Xiong, C., Benzinger, T.L., et al., 2012. Clinical and biomarker changes in dominantly inherited Alzheimer's disease. *N. Engl. J. Med.* 367 (9), 795–804.
- Bilgel, M., Jedynak, B., Wong, D.F., Resnick, S.M., Prince, J.L., 2015. Temporal trajectory and progression score estimation from voxelwise longitudinal imaging measures: application to amyloid imaging. In: *Proceedings of Information Processing in Medical Imaging*. Springer, pp. 424–436.
- Brookmeyer, R., Johnson, E., Ziegler-Graham, K., Arrighi, H.M., 2007. Forecasting the global burden of Alzheimer's disease. *Alzheimer's Dementia* 3 (3), 186–191.
- Cash, D.M., Frost, C., Theme, L.O., Ünay, D., Kandemir, M., Fripp, J., Salvado, O., Bourgeat, P., Reuter, M., Fischl, B., et al., 2015. Assessing atrophy measurement techniques in dementia: results from the MIRIAD atrophy challenge. *NeuroImage* 123, 149–164.
- Donohue, M.C., Jacqmin-Gadda, H., Le Goff, M., et al., 2014. Estimating long-term multivariate progression from short-term data. *Alzheimer's Dementia* 10 (5), S400–S410.
- Fontein, H.M., Clarkson, M.J., Modat, M., et al., 2011. An event-based disease progression model and its application to familial Alzheimer's disease. In: *IPMI*. Springer, pp. 748–759.
- Frisoni, G.B., Fox, N.C., Jack, C.R., Scheltens, P., Thompson, P.M., 2010. The clinical use of structural MRI in Alzheimer's disease. *Nat. Rev. Neurol.* 6 (2), 67–77.
- Gordon, B.A., Blazey, T., Su, Y., Fagan, A.M., Holtzman, D.M., Morris, J.C., Benzinger, T.L., 2016. Longitudinal  $\beta$ -amyloid deposition and hippocampal volume in preclinical Alzheimer's disease and suspected non-Alzheimer disease pathophysiology. *Jama Neurol.* 73 (10), 1192–1200.
- Guerrero, R., Schmidt-Richberg, A., Ledig, C., Tong, T., Wolz, R., Rueckert, D., 2016. Instantiated mixed effects modeling of Alzheimer's disease markers. *NeuroImage* 142, 113–125.
- Henneman, W., Sluiter, J., Barnes, J., Van Der Flier, W., Sluiter, I., Fox, N., Scheltens, P., Vrenken, H., Barkhof, F., 2009. Hippocampal atrophy rates in Alzheimer's disease. added value over whole brain volume measures. *Neurology* 72 (11), 999–1007.
- Hill, D.L., Schwarz, A.J., Isaac, M., Pani, L., Vamvakas, S., Hemmings, R., Carrillo, M.C., Yu, P., Sun, J., Beckett, L., et al., 2014. Coalition against major diseases/european medicines agency biomarker qualification of hippocampal volume for enrichment of clinical trials in pre-dementia stages of Alzheimer's disease. *Alzheimer's Dementia* 10 (4), 421–429.
- Jack, C.R., Knopman, D.S., Jagust, W.J., Shaw, L.M., Aisen, P.S., Weiner, M.W., Petersen, R.C., Trojanowski, J.Q., 2010. Hypothetical model of dynamic biomarkers of the Alzheimer's pathological cascade. *Lancet Neurology* 9 (1), 119–128.
- Jack Jr., C.R., Knopman, D.S., Chételat, G., Dickson, D., Fagan, A.M., Frisoni, G.B., Jagust, W., Mormino, E.C., Petersen, R.C., Sperling, R.A., et al., 2016. Suspected non-Alzheimer disease pathophysiology – concept and controversy. *Nat. Rev. Neurol.* 12, 117–124.
- Jedynak, B.M., Lang, A., Liu, B., Katz, E., Zhang, Y., Wyman, B.T., Raunig, D., Jedynak, C.P., Caffo, B., Prince, J.L., et al., 2012. A computational neurodegenerative disease progression score: method and results with the Alzheimer's disease neuroimaging initiative cohort. *NeuroImage* 63 (3), 1478–1486.
- Joensuu, H., Vehtari, A., Riihimäki, J., Nishida, T., Steigen, S.E., Brabec, P., Plank, L., Nilsson, B., Cirilli, C., Braconi, C., et al., 2012. Risk of recurrence of gastrointestinal stromal tumour after surgery: an analysis of pooled population-based cohorts. *Lancet Oncol.* 13 (3), 265–274.
- Kneip, A., Gasser, T., 1988. Convergence and consistency results for self-modeling nonlinear regression. *Ann. Statistics* 16, 82–112.
- Lázaro-Gredilla, M., Van Vaerenbergh, S., Lawrence, N.D., 2012. Overlapping mixtures of Gaussian processes for the data association problem. *Pattern Recognit.* 45 (4), 1386–1395.
- Liu, Q., Lin, K.K., Andersen, B., Smyth, P., Ihler, A., 2010. Estimating replicate time shifts using Gaussian process regression. *Bioinformatics* 26 (6), 770–776.
- Lorenzi, M., Donohue, M., Paternico, D., Scarpa, C., Ostrowitzki, S., Blin, O., Irving, E., Frisoni, G., 2010. Enrichment through biomarkers in clinical trials of Alzheimer's drugs in patients with mild cognitive impairment. *Neurobiol. Aging* 31 (8), 1443–1451.
- Lorenzi, M., Ziegler, G., Alexander, D.C., Ourselin, S., 2015. Efficient Gaussian process-based modelling and prediction of image time series. In: *IPMI*. Springer, pp. 626–637.
- Lorenzi, M., Simpson, I.J., Mendelson, A.F., Vos, S.B., Cardoso, M.J., Modat, M., Schott, J.M., Ourselin, S., 2016. Multimodal image analysis in Alzheimer's disease via statistical modelling of non-local intensity correlations. *Sci. Rep.* 6, 22161.
- Marinescu, R.V., Eshghi, A., Lorenzi, M., et al., 2017. A vertex clustering model for disease progression: application to cortical thickness images. In: *IPMI*. Springer (p. to appear).
- McKhann, G., Drachman, D., Folstein, M., Katzman, R., Price, D., Stadlan, E.M., 1984. Clinical diagnosis of Alzheimer's disease report of the nincds-adrcd work group\* under the auspices of department of health and human services task force on Alzheimer's disease. *Neurology* 34 (7), 939–939.
- Mendelson, A.F., Zuluaga, M.A., Lorenzi, M., Hutton, B.F., Ourselin, S., 2016. Selection bias in the reported performances of AD classification pipelines. *NeuroImage Clin.* 14, 400–416.
- Mormino, E.C., Papp, K.V., Rentz, D.M., Schultz, A.P., LaPoint, M., Amariglio, R., Hanseuw, B., Marshall, G.A., Hedden, T., Johnson, K.A., et al., 2016. Heterogeneity in suspected non-Alzheimer disease pathophysiology among clinically normal older individuals. *Jama Neurol.* 73 (10), 1185–1191.

- Mwangi, B., Tian, T.S., Soares, J.C., 2014. A review of feature reduction techniques in neuroimaging. *Neuroinformatics* 12 (2), 229–244.
- Nickisch, H., Rasmussen, C.E., 2008. Approximations for binary Gaussian process classification. *J. Mach. Learn. Res.* 9 (Oct), 2035–2078.
- R Core Team, 2015. *R: a Language and Environment for Statistical Computing*, R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>.
- Rasmussen, C.E., 2006. *Gaussian Processes for Machine Learning*. Springer.
- Riihimäki, J., Vehtari, A., 2010. Gaussian processes with monotonicity information. In: *AISTATS*, vol. 9, pp. 645–652.
- Ross, J.C., Dy, J.G., 2013. Nonparametric mixture of Gaussian processes with constraints. In: *Proceedings of International Conference of Machine Learning (ICML)*, pp. 1346–1354.
- Schiratti, J.-B., Allasonniere, S., Routier, A., Colliot, O., Durrleman, S., 2015. A mixed-effects model with time reparametrization for longitudinal univariate manifold-valued data. In: *IPMI*. Springer, pp. 564–575.
- Schmidt-Richberg, A., Guerrero, R., Ledig, C., Molina-Abril, H., Frangi, A.F., Rueckert, D., 2015. Multi-stage biomarker models for progression estimation in Alzheimer's disease. In: *International Conference on Information Processing in Medical Imaging*. Springer, pp. 387–398.
- Shinkins, B., Perera, R., 2013. Diagnostic uncertainty: dichotomies are not the answer. *Br. J. General Pract.* 63, 122–123.
- Wolz, R., Heckemann, R.A., Aljabar, P., Hajnal, J.V., Hammers, A., Lötjönen, J., Rueckert, D., 2010. Measurement of hippocampal atrophy using 4D graph-cut segmentation: application to ADNI. *NeuroImage* 52 (1), 109–118.
- Xu, Z., Shen, X., Pan, W., 2014. Longitudinal analysis is more powerful than cross-sectional analysis in detecting genetic association with neuroimaging phenotypes. *PLoS One* 9 (8), e102312.
- Yang, E., Farnum, M., Lobanov, V., et al., 2011. Quantifying the pathophysiological timeline of Alzheimer's disease. *J. Alzheimer's Dis.* 26 (4), 745–753.
- Younes, L., Albert, M., Miller, M.I., Team, B.R., et al., 2014. Inferring changepoint times of medial temporal lobe morphometric change in preclinical Alzheimer's disease. *NeuroImage Clin.* 5, 178–187.
- Young, J., Modat, M., Cardoso, M.J., Mendelson, A., Cash, D., Ourselin, S., 2013. Accurate multimodal probabilistic prediction of conversion to Alzheimer's disease in patients with mild cognitive impairment. *NeuroImage Clin.* 2, 735–745.
- Young, A.L., Oxtoby, N.P., Daga, P., Cash, D.M., Fox, N.C., Ourselin, S., Schott, J.M., Alexander, D.C., 2014. A data-driven model of biomarker changes in sporadic Alzheimer's disease. *Brain* 137 (9), 2564–2577.
- Yu, P., Sun, J., Wolz, R., Stephenson, D., Brewer, J., Fox, N.C., Cole, P.E., Jack, C.R., Hill, D.L., Schwarz, A.J., et al., 2014. Operationalizing hippocampal volume as an enrichment biomarker for amnesic mild cognitive impairment trials: effect of algorithm, test-retest variability, and cut point on trial cost, duration, and sample size. *Neurobiol. Aging* 35 (4), 808–818.