# S-EBM: Generalising event-based modelling of disease progression for simultaneous events

Parker CS[1,⊥], Oxtoby NP[1], Alexander DC[1], Zhang H[1] for the Alzheimer's Disease Neuroimaging Initiative[2]

[1] Centre for Medical Image Computing, Department of Computer Science, UCL, London, UK

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

[⊥] Corresponding author: christopher.parker@ucl.ac.uk. University College London, Centre for Medical Image Computing, 90 High Holborn, Floor 1, London, UK, WC1V6LJ

1

## Abstract

Estimating the temporal evolution of biomarker abnormalities in disease informs understanding of early disease processes and facilitates subject staging, which may augment the development of early therapeutic interventions and provide personalised treatment tools. Event-based modelling of disease progression (EBM) is a data-driven technique for inferring a sequence of biomarker abnormalities, or events, from cross-sectional or short-term longitudinal datasets and has been applied to a variety of different diseases, including Alzheimer's disease. Conventional EBM (C-EBM) assumes the sequence of biomarker abnormalities occurs in series, with one biomarker event per disease progression stage. However, events may occur simultaneously, for example due to the presence of shared causal factors, a property which cannot be inferred from C-EBM. Here we introduce simultaneous EBM (S-EBM), a generalisation of C-EBM to enable estimation of simultaneous events. S-EBM can estimate a wider range of sequence types than C-EBM while being fully backward compatible with the original model. Using simulated data, we firstly demonstrate the inability of C-EBM to infer simultaneous events. We next assess the accuracy of S-EBM against ground truth data and subsequently demonstrate a real-world example application to sequence disease progression in Alzheimer's disease. Simulations show that C-EBM can not discern serial events with high biomarker variance from simultaneous events, preventing its use for inferring simultaneous events. S-EBM has high estimation accuracy against ground truth for a range of sequence types (fully simultaneous, partially simultaneous, serial), number of biomarkers and biomarker variances. When applied to Alzheimer's disease biomarker data from ADNI, S-EBM estimated a sequence where events within sets of biomarker domains occur simultaneously. Accumulation of total and phosphorylated tau in cerebrospinal fluid; performance on RAVLT, ADAS-Cog and MMSE cognitive test scores; and volumetric decline in temporal regional brain volumes, were better described as groups of simultaneous events rather than a single set of serial events (likelihood ratio >> 1,000). Furthermore, C-EBM may be confidently incorrect regarding the serial ordering. S-EBM may be applied to prospective and retrospective biomarker data to refine understanding of disease progression and generate new hypotheses regarding disease aetiology and spread.

## 1. Introduction

Estimating the temporal progression of biomarker abnormalities throughout the course of a disease identifies biomarkers of early disease, which generates hypotheses regarding disease aetiology and spread; and facilitates subject staging, which may aid the development of therapeutic interventions and personalised treatment.

Disease progression has been previously estimated using hypothesis-driven approaches following literature review or post-mortem examination. For example, the Jack curves (Jack Jr. et al 2010) describe the evolution of biomarkers abnormalities in Alzheimer's disease (AD), and Braak stages were derived from post-mortem examination of AD patients (Braak, H & Braak, E 1991). Although these approaches are informative, they are qualitative in nature. Data-driven approaches are needed for objective assessment of disease spread. In the ideal scenario the temporal trajectory of different biomarkers is derived from longitudinal data acquired throughout the disease course. However, in practice cross-sectional or short-term longitudinal data are the predominant type of biomarker data available. There is therefore a need for approaches that estimate disease progression from such data.

Event-based modelling of disease progression (EBM) is a data-driven approach that estimates the evolution of biomarker abnormalities from cross-sectional or short-term longitudinal data (Fonteijn et al 2012). EBM has been applied to estimate progression of biomarker abnormality in a variety of diseases, including AD (Young et al 2014), Huntington's disease (Wijeratne et al 2018), multiple sclerosis (Eshaghi et al 2018) and amyotrophic lateral sclerosis (Gabel et al 2020).

Underlying the conventional EBM (C-EBM) approach (Fonteijn et al 2012), as well as its recent variants, is the assumption that biomarker abnormalities are ordered serially, i.e. no two biomarkers may become abnormal concurrently. However, biomarker abnormalities may occur simultaneously when they are driven by common causative factors, or be better approximated as simultaneous than as serial when the difference between their temporal trajectories is unresolvably small. Such simultaneous events cannot be inferred from C-EBM as they are excluded from the model by construction. The positional uncertainty that C-EBM estimates may suggest the presence of simultaneous events, but can also simply reflect high variance in biomarker measurements. By not accounting for simultaneous events, C-EBM may incorrectly estimate the sequence and patient staging, limiting its ability to impact disease understanding and therapeutic development.

To overcome this limitation, we introduce simultaneous EBM (S-EBM), a generalisation of C-EBM that can estimate a sequence containing simultaneous events. By allowing simultaneous events, a wider range of disease progression models can be

107 estimated from any given biomarker data input. In this study, we demonstrate C-EBM's

108 inability to infer simultaneous events, describe the theory of S-EBM and sequence

109 estimation, evaluate the performance of S-EBM against ground truth synthetic data, and

110 provide an example application to sequence evolution of biomarker abnormalities in AD. We

111 show that S-EBM can reliably estimate sequences containing simultaneous events and that

112 such a sequence can better explain the evolution of AD biomarker abnormality.

113

## 2. Theory

115 2.1. Generalising the event-based model

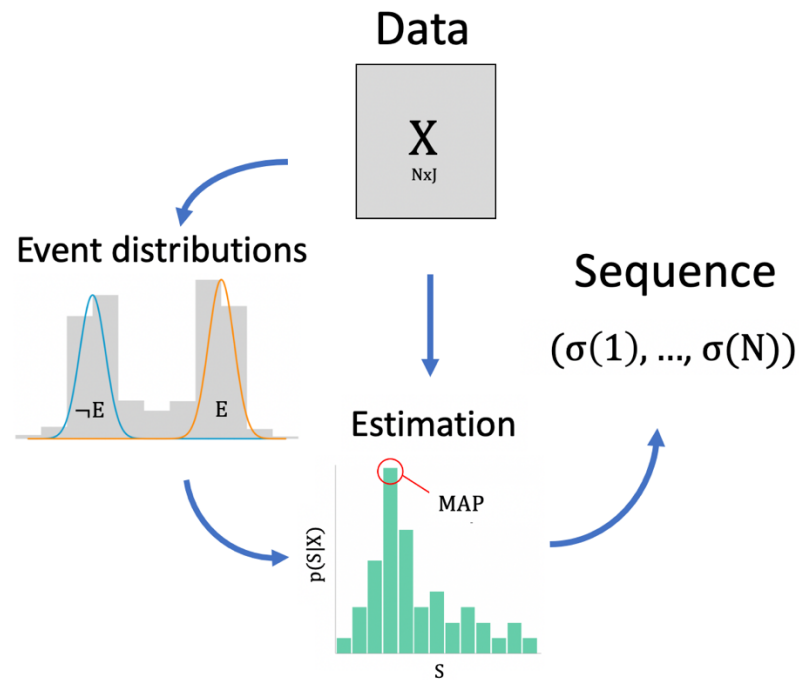116 *2.1.1. Overview of the conventional event-based model*

117       C-EBM represents the progression of biomarker abnormalities in disease by a

118 sequence, which is an ordered list that encodes the temporal order in which each biomarker

119 undergoes a transition from a normal state to an abnormal state. These transitions, termed

120 events, demarcate the disease progression stages, from which subjects are assumed to be

121 uniformly sampled.

122       A key assumption of C-EBM is monotonicity of biomarker evolution i.e. that

123 biomarkers transition to an abnormal state but do not subsequently revert. Thus, in the first

124 stage all biomarkers are in a normal state and at each subsequent stage a biomarker

125 transitions to an abnormal state, until the final stage where all biomarkers are abnormal. A

126 further key assumption of C-EBM is that all subjects are sampled from the same disease

127 trajectory. In other words, the set of biomarker measurements for a given subject provides a

128 snapshot of the disease at a particular stage. Furthermore, the subjects are assumed to be

129 sampled from a single disease progression sequence.

130       C-EBM seeks the sequence with highest posterior probability given the observed

131 biomarker measurements. By assuming an equal prior probability for all possible sequences,

132 this becomes equivalent to the sequence likelihood i.e. the probability of the data given the

133 sequence. As the sequence prescribes the set of events for each disease stage, then given

134 the probability density functions associated with each biomarkers' possible event state (see

135 section 2.4. Event distributions), then the likelihood of the sequence can be evaluated and

136 subsequently maximised across sequence samples. A summary of sequence estimation is

137 shown in Fig. 1.

138

139

**Figure 1.** Overview of sequence estimation in C-EBM. C-EBM finds the sequence, $S$, with maximum posterior probability given the biomarker measurements, $X$. Given an equal prior probability of each sequence, this is equivalent to the maximum likelihood sequence.

The sequence likelihood is equal to the joint probability of observing the set of subjects' data. Given the input data matrix $X$, an N-by-J matrix containing N biomarker measurements for J subjects, and assuming that each subject is sampled independently, then the likelihood of the sequence, $S$, is the product of subject probabilities:

$$p(X|S) = \prod_{j=1}^{J} p(X_j|S) \tag{1}$$

where $X_j$ is a column of $X$ corresponding to the N biomarker measurements for subject j.

As described below, the formulation of $p(X_j|S)$ makes reference to the set of event states' distributions at each disease stage. As these events are defined by $S$, the formulation of $p(X_j|S)$ depends on the specific form of $S$. Next, we describe how the sequence is specified and likelihood formulation derived for C-EBM, which assumes the events occur in series, before describing the generalisation of the sequence and likelihood formulation for simultaneous events.

*2.1.2. Conventional event-based model: sequence specification and likelihood function*

C-EBM specifies the sequence as a permutation of the biomarker indices $1, ..., N$. Each element of $S$, $s(i)$, holds the biomarker event occurring at the i'th disease progression stage. For example, for a sequence of four biomarkers a possible sequence is $S = (2,3,4,1)$,

5

162  which describes a disease progression where the first biomarker abnormality occurs in

163  biomarker 2, followed by biomarker 3, then biomarker 4 and finally biomarker 1.

164  With each biomarkers' event states written as $\neg E$ for normal and $E$ for abnormal, then

165  at a particular stage, k, of the sequence the events have occurred for $E_{s(i)}, \ldots, E_{s(k)}$ but have

166  not yet occurred for $E_{s(k+1)}, \ldots, E_{s(N)}$. Given independence of biomarker measurements for

167  the combination of events at each sequence position, the subjects' probability given the

168  sequence and stage, k, is written as:

169

170
$$p(X_j|S, k) = \prod_{i=1}^{k} p\big(x_{s(i),j} \mid E_{s(i)}\big) \prod_{i=k+1}^{N} p\big(x_{s(i),j} \mid \neg E_{s(i)}\big) \qquad (2)$$

171

172  Because each subjects' position in the sequence is considered unknown a priori, it is

173  marginalised out over each possible position:

174
$$p(X_j|S) = \sum_{k=0}^{N} p(k)p(X_j|S, k) \qquad (3)$$

175

176  The prior probability of each position, $p(k)$, is assumed to be constant and defined as

177  $\frac{1}{N+1}$, where $N + 1$ (or equivalently $|S| + 1$) is the number of stages. By substituting Eq. 2 into

178  3, then the total likelihood defined in Eq. 1 is written as:

179

180
$$p(X|S) = \prod_{j=1}^{J} \left( \sum_{k=0}^{N} p(k) \left[ \prod_{i=1}^{k} p\big(x_{s(i),j} \mid E_{s(i)}\big) \prod_{i=k+1}^{N} p\big(x_{s(i),j} \mid \neg E_{s(i)}\big) \right] \right) \qquad (4)$$

181

182  Because the sequence can contain only one biomarker event at each position, it

183  cannot represent simultaneous events.

184

185  *2.1.3. Simultaneous event-based model: sequence specification and likelihood function*

186  To generalise C-EBM for simultaneous events, the sequence specification is updated

187  from an ordered list of biomarker indices to an ordered list of sets. Each set, $s(i)$, contains

188  one or more biomarker indices corresponding to the events at position i in the sequence. For

189  example, for four biomarkers a sequence containing only serial events is written $S =$

190  $(\{2\}, \{1\}, \{3\}, \{4\})$ and a sequence containing simultaneous events is written $S =$

191  $(\{2\}, \{1,3\}, \{4\})$. Given the length of the sequence can vary, the number of positions in the

192  sequence is now defined as $|S| + 1$ instead of $N + 1$. Therefore, the prior probability of each

6

position in the sequence is $p(k;\ S) = \frac{1}{|S|+1}$ and the likelihood of each subjects' data given their position is unknown a priori is written as:

$$p(X_j|S) = \sum_{k=0}^{|S|} p(k;\ S)p(X_j|S, k) \tag{5}$$

As before, the likelihood of each subjects' data given their position $k$, $p(X_j|S, k)$ is the joint probability over the subjects' biomarker values given each biomarkers event state at that sequence position. For a position $k$ in the sequence, the events have occurred for biomarkers $\cup_{1\leq i\leq k} s(i)$, whereas the events have not occurred for biomarkers $\cup_{k<i\leq|S_m|} s(i)$. Hence, the likelihood of each subjects' data given their position is written as:

$$p(X_j|S, k) = \prod_{\substack{m\in \\ \cup_{1\leq i\leq k} s(i)}} p(x_{m,j} \mid E_m) \prod_{\substack{m\in \\ \cup_{k<i\leq|S|} s(i)}} p(x_{m,j} \mid \neg E_m) \tag{6}$$

By substituting Eq. 6 into 5, then the total likelihood defined by Eq. 1 is written as:

$$p(X|S) = \prod_{j=1}^{J} \left( \sum_{k=0}^{|S|} p(k; S) \left[ \prod_{\substack{m\in \\ \cup_{1\leq i\leq k} s(i)}} p(x_{m,j}|E_m) \prod_{\substack{m\in \\ \cup_{k<i\leq|S|} s(i)}} p(x_{m,j}|\neg E_m) \right] \right) \tag{7}$$

This likelihood formulation is a fully generalised form of the C-EBM but can represent a wider range of sequence types. In the case of serial events, the likelihood defined in Eq. 7 becomes equal to the C-EBM likelihood defined in Eq. 4.

## 2.2. Sequence estimation

### 2.2.1. Conventional event-based model

In C-EBM (Fonteijn et al 2012), the sequence is estimated as the characteristic ordering of biomarker events, which is the average position of each event following Markov Chain Monte Carlo (MCMC) sampling of $p(S|X)$. In subsequent work (Young et al 2014), a stochastic greedy ascent was used to estimate the maximum likelihood sequence. As we aimed to compare the sequence obtained from (Young et al 2014) between C-EBM and S-EBM, this is the approach we adopt here.

223    The greedy ascent proceeds by iteratively perturbing the sequence and retaining

224    those with higher likelihood for some given number of iterations. At each iteration, a

225    perturbation of the sequence is generated by swapping the positions of two biomarker

226    events. For example, the if the current sequence is $(2, 3, 4, 1)$, then a perturbed sequence

227    can be generated by swapping biomarkers 4 and 2, giving the sequence $(4, 3, 2, 1)$. To

228    prevent dependence of the greedy ascent on the initial random sequence, a number of

229    initialisations are performed and the sequence with maximum likelihood over all ascents is

230    the estimated sequence.

231

232    *2.2.2. Simultaneous event-based model*

233    To enable traversal of the full space of sequences that contain any combination of

234    simultaneous events, we update the sequence perturbation method: a biomarker is chosen

235    at random and is replaced at any other valid position in the sequence. For example, if the

236    sequence is $(\{2\}, \{1, 3\}, \{4\})$, then a perturbed sequence can be generated by randomly

237    choosing biomarker 3 and replacing the biomarker at position 4 in the sequence, giving the

238    sequence $(\{2\}, \{1\}, \{4\}, \{3\})$. Other possible perturbations are shown in Supplementary Table

239    1. This perturbation method is compatible with the MCMC sampling method described in

240    (Fonteijn 2012, Young et al 2014), as it retains the property of symmetric transition

241    probability $p(S_{t+1}|S_t) = p(S_t|S_{t+1})$, which simplifies the formulation of the acceptance

242    probability.

243

244    2.3. Event state probability density functions

245    Calculating the sequence likelihood requires the probability density functions of each

246    biomarker under the condition that the event has or has not occurred,

247    $p(x_{m=1,j}|E_{m=1}), \dots, p(x_{m=N,j}|E_{m=N})$ and $p(x_{m=1,j}|\neg E_{m=1}), \dots, p(x_{m=N,j}|\neg E_{m=N})$, respectively.

248    Hypothetically, if each subjects' position in the sequence is known, then the event

249    state for each biomarker measurement is also known. For example, for a given biomarker $i$

250    and its event state $E_i$, the probability density function $p(x_{i,j}|E_i)$ can be fitted to the

251    measurements $\{x_{i,j} \mid k(j) \geq p, s(p) = i\}$ (i.e. the measurements for the subjects at a position

252    greater or equal to the position of the event for biomarker $i$), where $k(j)$ is the position in the

253    sequence of subject $j$.

254    However, as the subjects' sequence position is unknown a priori, then the

255    assumption is made that the measurements are drawn from a mixture distribution $p(x_{i,j}) =$

256    $w_i p(x_{i,j}|E_i) + (1 - w_i)p(x_{i,j}|\neg E_i)$, whose components are then recovered by fitting a mixture

257    model to all measurements $\{x_{i,j}|j = 1, \dots N\}$.

258

8

## 3. Materials and Methods

3.1. Simulation experiments

*3.1.1. Simultaneous event-based forward model*

A forward model is used in this study to generate biomarker data for simulation experiments. The model generates data from a given ground truth sequence that can contain simultaneous events. The required inputs to the forward model are (i) the sequence (as described in 2.1.3. Simultaneous event-based model: sequence specification and likelihood function), (ii) the event distributions for each biomarker and (iii) the number of datapoints (i.e., subjects) to sample.

Firstly, a position $k$, of the subject within the disease progression sequence is sampled from the uniform prior distribution $\mathrm{Unif}\{0, |S|\}$. The biomarker data for this subject, indexed by $j$, is then generated by sampling from the event distributions corresponding to the position in the sequence: $\sim x_{m,j}|E_m$ if $m \in \bigcup_{1 \le i \le k} s(i)$ , or $\sim x_{m,j}|\neg E_m$ if $m \in \bigcup_{k \le i \le |S|} s(i)$. The process is then repeated for the specified number of subjects, returning a matrix $X$ of size N-by-J containing the data samples for $J$ subjects and $N$ biomarkers.

*3.1.2 Experiment 1: biomarker variance, simultaneous events and C-EBM uncertainty*

To demonstrate that the uncertainty in event positions derived from C-EBM cannot be used to infer the presence of simultaneous events, we quantified the effect of both biomarker variance and simultaneous events on degree of sequence uncertainty. We hypothesised that both biomarker variance and simultaneous events can separately result in a high degree of uncertainty in event positions.

Data was simulated for two biomarkers sampled from either a serial event sequence ($\{1\}, \{2\}$), or simultaneous event sequence ($\{1,2\}$), whose probability density functions were gaussian with a mean of zero for the normal event states (Eqs. 8 and 9) and one for abnormal event states (Eqs. 10 and 11). Standard deviation was varied from 0.05 to 2.00 and was equal for each biomarker and event state.

$$p(x_{1,j}|\neg E_1) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{(x_{1,j}-0)^2}{4\pi^2}\right) \tag{8}$$

$$p(x_{2,j}|\neg E_2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{(x_{2,j}-0)^2}{4\pi^2}\right) \tag{9}$$

$$p(x_{1,j}|E_1) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{(x_{1,j}-1)^2}{4\pi^2}\right) \tag{10}$$

$$p\left(x_{2,j}\middle|E_2\right) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{\left(x_{2,j} - 1\right)^2}{4\pi^2}\right) \tag{11}$$

For each sequence and standard deviation combination, one hundred datasets were simulated, each with ten 'control' subjects at position zero, where no events have yet occurred, ten 'end-stage patients' at the final sequence position $|S|$, where all events have occurred, and twenty 'intermediate-stage patients', which are sampled uniformly from the sequence positions i.e., $k \sim \mathrm{Unif}\{0, |S|\}$. To remove the added variability in positional uncertainty due to the estimation of event distributions, these distributions were determined from their simulation definitions.

For each simulated dataset, the uncertainty was quantified in a positional variance matrix, $P$, whose $i, j'$th entry gives the probability that biomarker $i$ is at position $j$. This probability is defined as the frequency over MCMC samples where biomarker $i$ is at position $j$ (Fonteijn et al 2012) i.e. $P_{i,j} = (\sum_{S \in S_{ij}} 1)/N_{mcmc}$, where $N_{mcmc}$ is the number of MCMC samples and $S_{ij}$ is the set of sequences where biomarker $i$ is at position $j$. In the case of a serial sequence containing only two biomarkers, this simplifies to $P_{i,j} = P(X|S_{i@j})$, where $S_{i@j}$ refers to the sequence with biomarker $i$ at position $j$. A binary decision was then made as to whether each positional variance matrix has a significant level of uncertainty or not. A significant level of uncertainty was defined as the highest probability in the matrix being less than 0.95, which corresponds to the absence of certainty (with 0.95 probability of higher) in biomarker positions. The proportion of matrices containing significant levels of uncertainty for the serial or simultaneous sequences was then plotted as a function of biomarker standard deviation.

*3.1.3. Experiment 2: Evaluation of simultaneous EBM performance*

We evaluated simultaneous EBM performance against a known ground truth sequence by quantifying the percentage of correctly estimated sequences over a set of one hundred simulations of biomarker data. The set of one hundred simulations was repeated for each combination of sequence type (serial, partially simultaneous and fully simultaneous), number of biomarkers (2, 4 and 10), number of subjects (40, 80 and 160) and biomarker variance (s.d's of 0.1, 0.2 and 0.3).

For each number of subjects, the subject types were split in a 1:2:1 ratio between control, intermediate and end-stage. As in Experiment 1 (section 3.1.2.), the means of the event states used to generate the simulated data were zero and one for normal and abnormal event states, respectively, and the standard deviations were equal for the biomarker event states for each s.d. value.

10

325    To sufficiently sample the set of possible sequences during sequence estimation, the

326    number of initialisations and iterations of the greedy ascent was adjusted for each number of

327    biomarkers: 1 and 2 respectively for two biomarkers, 10 and 100 respectively for 4

328    biomarkers, and 50 and 1000 respectively for 10 biomarkers. For all sequence estimations,

329    the event distributions were fitted using the data from the control and end-stage subjects.

330

### 3.1.4. Experiment 3: Comparison to conventional EBM for serial events

332    To evaluate the ability of S-EBM to correctly identify a sequence containing serial

333    events in the case where C-EBM reports high uncertainty, we quantify the percentage of

334    correctly estimated sequences as a function of biomarker variance for the range of

335    biomarker variance that resulted in a high proportion of positional uncertainty, as determined

336    from section 3.1.2. The simulation conditions are as described in 3.1.2. except with the

337    sequence estimation being performed by either C-EBM or S-EBM on the serial sequence.

338

## 3.2. Application to Alzheimer's disease progression

340    We applied S-EBM to sequence the evolution of biomarker abnormalities in AD while

341    accounting for simultaneous events and compared it to the serial sequence estimated by C-

342    EBM. Our pipeline for data selection follows that of (Young et al 2014) but utilises existing

343    sources of pre-compiled AD data.

344

### 3.2.1. AD biomarker source

346    Biomarkers of cerebrospinal fluid (CSF) (total tau, phosphorylated tau, amyloid-$\beta_{1\text{-}42}$),

347    cognitive test scores (RAVLT, ADAS-Cog, MMSE) and regional brain volumes

348    (hippocampus, entorhinal cortex, mid-temporal gyrus, fusiform and ventricles) were obtained

349    from the TADPOLE dataset, which is available for download from the Alzheimer's disease

350    Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu) (Mueller et al 2005). TADPOLE

351    is a pre-compiled source of ADNI biomarker data that includes data from phases 1, GO and

352    2 of ADNI. TADPOLE datasets D1 and D2, which contain biomarker data from every

353    individual that has participated in in at least two separate visits, were used in this study. The

354    image processing steps used by ADNI to generate the biomarkers later compiled in the

355    TADPOLE dataset are described in 3.2.2. ADNI processing.

356

### 3.2.2. ADNI processing

358    CSF measurements of total tau, phosphorylated tau and amyloid-$\beta$ were obtained via

359    lumbar puncture (Shaw et al 2009). Cognitive test scores were obtained via specialist clinical

360    assessment (Crane et al 2012). Structural magnetic resonance (MR) images were acquired

11

and underwent pre-processing with standard ADNI pipelines (Jack Jr. et al 2008), which involved correction for gradient non-linearity, B1 non-uniformity correction and peak sharpening. Regional volumes were extracted using Freesurfer cross-sectional and longitudinal pipelines (Reuter et al 2012).

### 3.2.3. Biomarker processing

Following (Young et al 2014), we included subjects with available biomarker data acquired at baseline up to 5[th] February 2013 from those subjects scanned at 1.5T. Brain volumes were averaged over hemispheres and normalised by intracranial volume to control for individual differences in head size. CSF total tau and phosphorylated tau were log-transformed to improve event distribution estimation. Cognitively normal subjects who were positive for CSF amyloid-$\beta$ (<992 pg/ml) or phosphorylated tau (>25 pg/ml) were removed to improve the estimation of event distributions, which are presumed to be predominantly normal in this group.

### 3.2.4. Event distributions

For each biomarker, probability density functions corresponding to the event having occurred or having not occurred, were fitted to the cognitively normal and AD patients' biomarker data using a constrained gaussian mixture model implemented in MATLAB, as described in (Young et al 2014). The standard deviations of each event component ($E$ and $\neg E$) are constrained to be less than or equal to that of the cognitively normal or AD group, respectively, and the means are constrained to be no less extreme than the cognitively normal or AD groups. These constraints ensure a robust fit in the case where the distributions of healthy and patient population overlap significantly.

### 3.2.5. Sequencing allowing simultaneous events

The maximum likelihood S-EBM sequence was estimated from 1,000,000 MCMC samples. MCMC was initialised using the sequence estimated from a greedy ascent performed with 200 initialisations each with 2,000 iterations.

### 3.2.6. Sequencing of serial events

The maximum likelihood C-EBM sequence was estimated using greedy ascent with 200 random initialisations, each with 2,000 iterations. 1,000,000 MCMC samples were taken to estimate the uncertainty in each biomarkers position.

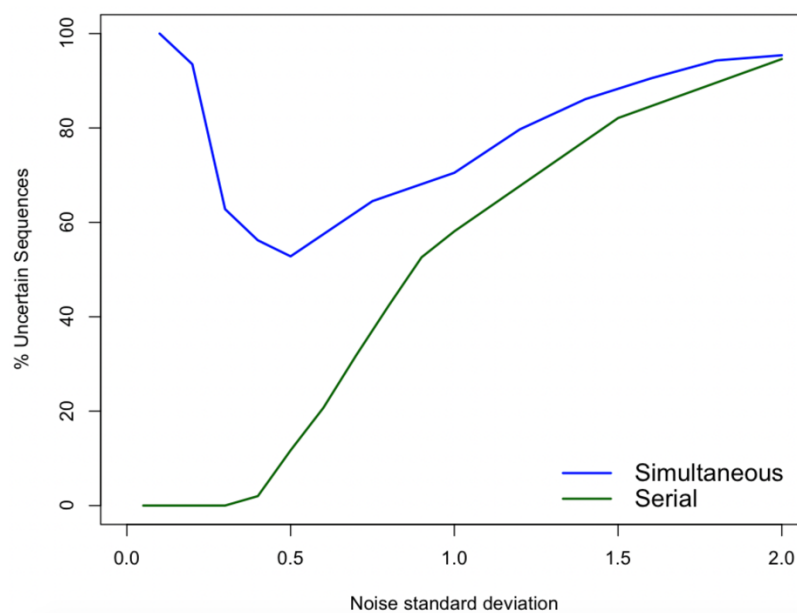# 4. Results and Discussion

4.1. Simulation experiments

*4.1.1. Experiment 1: biomarker variance, simultaneous events and C-EBM uncertainty*

A serial sequence with high biomarker variance can produce data which is interpreted by C-EBM as having high positional uncertainty (Fig. 2, green line). This uncertainty arises from the relative smoothness of the likelihood function across the sequence space due to overlapping event probability distributions. However, the same degree of uncertainty is also apparent in data produced from sequences containing simultaneous events (Fig. 2, blue line). This many-to-one mapping between sequence features (biomarker variance, simultaneous events) and positional uncertainty suggests that the presence of positional uncertainty in a particular dataset does not imply that the sequence contains simultaneous events. This prevents the use of C-EBM's positional uncertainty for detecting sequences containing simultaneous events.



**Figure 2.** The relation between biomarker standard deviation (x-axis) and uncertainty in the serial sequence estimated by C-EBM (y-axis) for simultaneous events (blue line) and serial events (green line). Both high biomarker variance in serial sequences, and sequences containing simultaneous events, result in a high percentage of uncertain sequences.
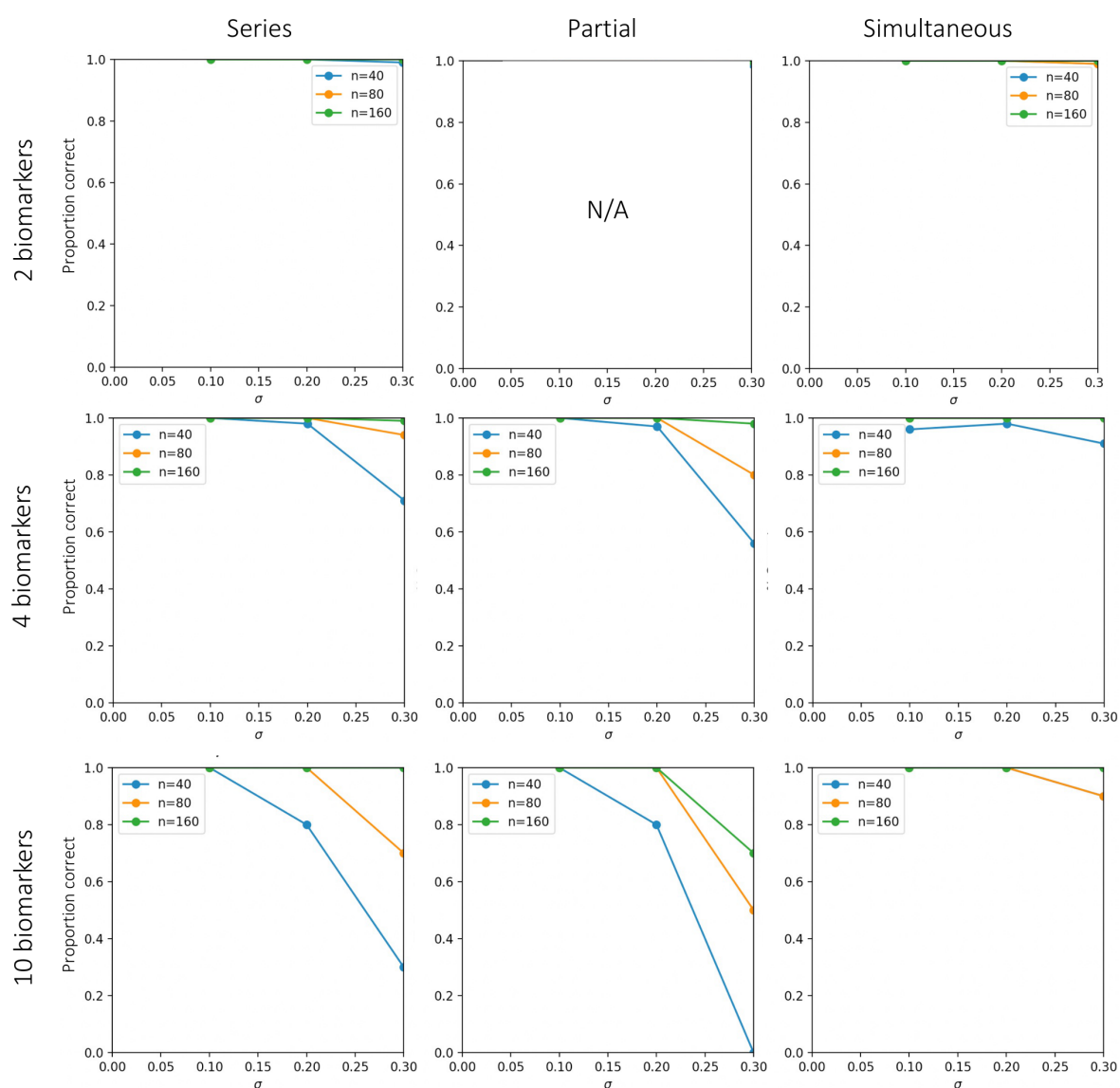
*4.1.2. Experiment 2: Evaluation of simultaneous EBM performance*

S-EBM accurately estimated sequences containing serial events, simultaneous events or both, under a range of experimental conditions (Fig. 3). Sequence estimation accuracy was high for sequences of 10 biomarkers and high biomarker variance when a sufficiently high number of datapoints was sampled. When fewer than 10 datapoints were

13

422 sampled per sequence position, accuracy tended to decrease for biomarker standard

423 deviations exceeding 0.1 for both serial and partially simultaneous sequences. Accuracy

424 was high for sequences containing simultaneous events under all conditions.

425      These results suggest that for moderately sized cohorts of individuals, S-EBM will

426 produce accurate estimates of sequences containing serial events, simultaneous events or

427 both. Given the increasing availability of large prospective and retrospective repositories of

428 cross-sectional or short-term longitudinal biomarker data, this technique has the potential to

429 inform on disease spread patterns for a range of disease. Of particular interest is using

430 retrospective data to provide a refined understanding of disease progression previously

431 estimated using C-EBM.

432



434 **Figure 3.** Accuracy of S-EBM sequence estimation for different sequence types (columns),

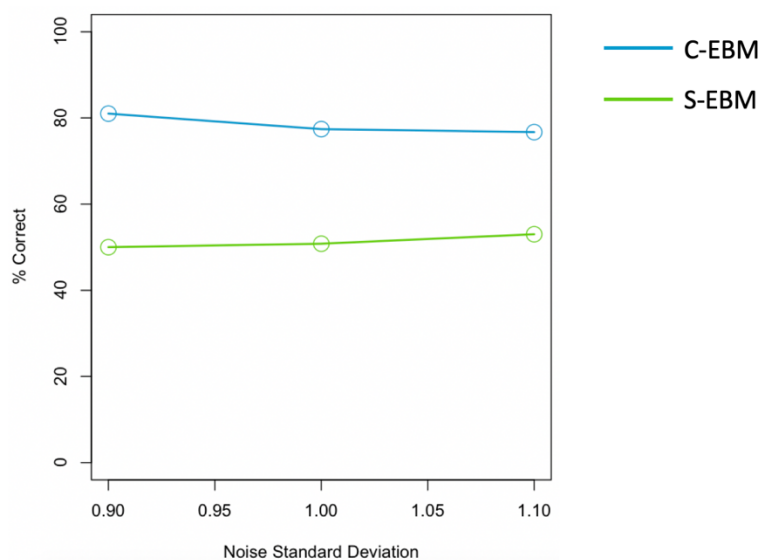435 numbers of biomarkers (rows), noise standard deviations (x-axis) and number of subjects

436    (coloured lines). Accuracy was high for almost all simulations but tended to decrease with

437    fewer subjects, higher noise standard deviation and more biomarkers.

438

439    *4.1.3. Experiment 3: Comparison to conventional EBM for serial events*

440        C-EBM had higher sequence estimation accuracy than S-EBM for noisy serial

441    sequences which had high C-EBM positional uncertainty (Fig. 4). This suggests that when

442    C-EBM is uncertain on the positional orderings, its maximum likelihood sequence is

443    nevertheless more likely to be correct than the maximum likelihood sequence estimated by

444    S-EBM. This may be expected given that the size of the sequence space of simultaneous

445    events is greater than that for serial sequences, which leaves more scope for false positives.

446    Despite this, without a priori knowledge of the sequence type, S-EBM offers the opportunity

447    to correctly identify a far wider range of types of sequences beyond those restricted by serial

448    order.

449



450

451    **Figure 4.** A comparison between C-EBM and S-EBM of serial sequence estimation

452    accuracy in the case where C-EBM reports high positional uncertainty. In this case C-EBM's

453    performance is superior to S-EBM due to the smaller sequence search space.
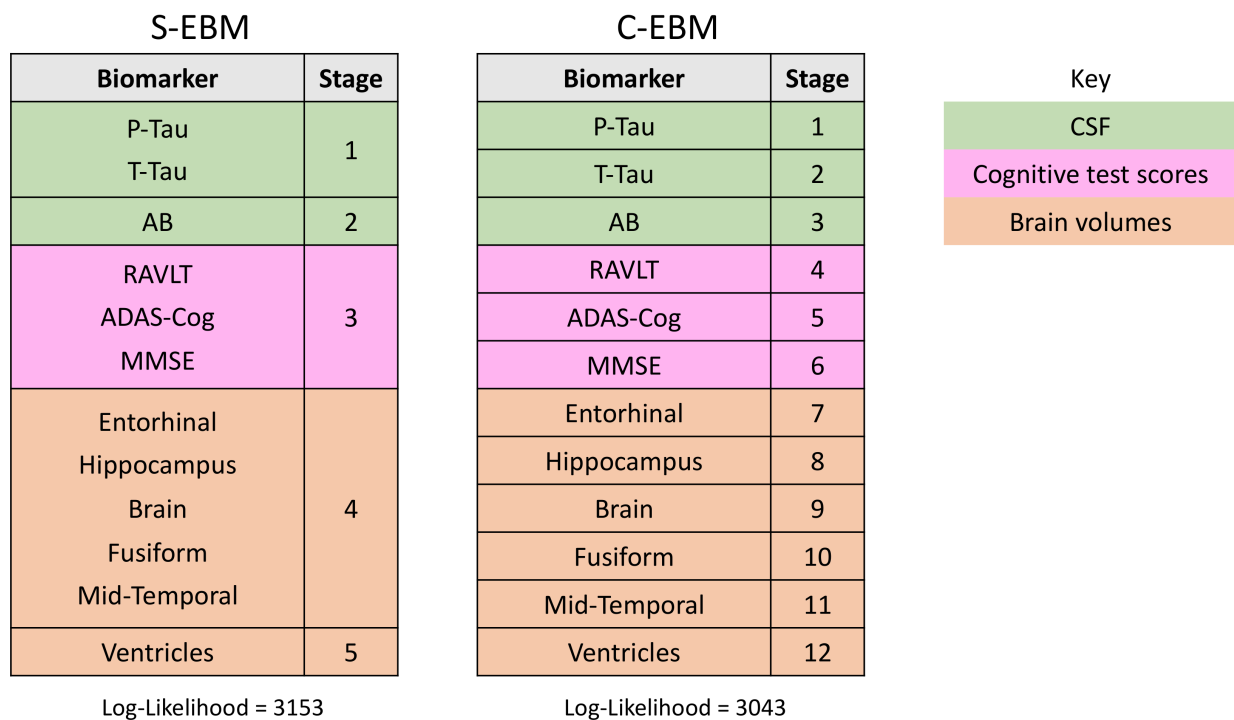
454

455

456

457

458

459

460

15

## 4.2. Application to Alzheimer's disease progression

### 4.2.1. S-EBM: estimated sequence allowing simultaneous events

The sequence of AD biomarker progression estimated by S-EBM is shown in Fig. 5. S-EBM identified a sequence containing simultaneous events which had a substantially higher log-likelihood compared to the serial sequence estimated by C-EBM.

Simultaneous events were estimated for biomarkers within common biomarker marker domains - CSF, cognitive test scores and brain volumes. Increased CSF total tau and phosphorylated tau were the first events in the sequence, occurring simultaneously, and were followed by high CSF amyloid-β. At disease stage three, low-scoring performance on cognitive test scores RAVLT, ADAS-Cog and MMSE were estimated as simultaneous events. Following cognitive events, the next disease stage consisted of simultaneous volumetric decline in temporal lobe brain regions. The final event in the sequence was increased ventricular volume.

### S-EBM

| Biomarker | Stage |
|---|---|
| P-Tau | 1 |
| T-Tau | |
| AB | 2 |
| RAVLT | 3 |
| ADAS-Cog | |
| MMSE | |
| Entorhinal | 4 |
| Hippocampus | |
| Brain | |
| Fusiform | |
| Mid-Temporal | |
| Ventricles | 5 |

Log-Likelihood = 3153

### C-EBM

| Biomarker | Stage |
|---|---|
| P-Tau | 1 |
| T-Tau | 2 |
| AB | 3 |
| RAVLT | 4 |
| ADAS-Cog | 5 |
| MMSE | 6 |
| Entorhinal | 7 |
| Hippocampus | 8 |
| Brain | 9 |
| Fusiform | 10 |
| Mid-Temporal | 11 |
| Ventricles | 12 |

Log-Likelihood = 3043

### Key

| | |
|---|---|
| CSF | |
| Cognitive test scores | |
| Brain volumes | |

**Figure 5.** The sequence of abnormality in biomarkers of CSF, cognitive test scores and brain volumes in AD estimated using S-EBM (left) and C-EBM (right). S-EBM estimates a sequence with substantially higher log-likelihood than C-EBM, by grouping certain biomarkers within domains into the same disease stage. In contrast, S-EBM assumes each biomarker abnormality occurs in series.

16

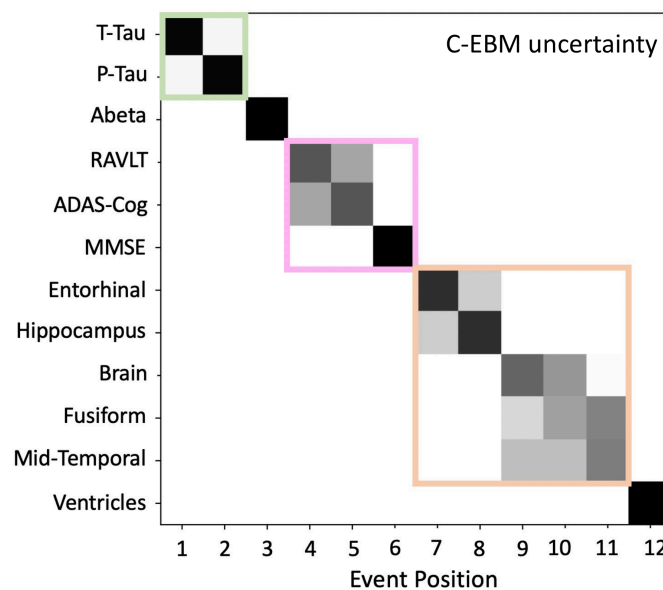### *4.2.2. C-EBM: estimated serial sequence*

The serial sequence estimated by C-EBM (Fig. 6) identified a lower log-likelihood sequence that, by design, assumed all events occur in series. However, it was consistent with the S-EBM sequence in finding a positional separation between groups of biomarker events belonging to different biomarker domains.

The C-EBM positional variance diagram (Fig. 6) however shows a heterogeneous distribution of positional uncertainty for the groups of simultaneous events, highlighting that positional uncertainty cannot be used to infer simultaneous events. Interstingly, interpreting the blocks of positional uncertainty as simultaneous events derives a sequence ({T-Tau}, {P-Tau}, {Abeta}, {RAVLT, ADAS-Cog}, {MMSE}, {Entorhinal, Hippocampus}, {Brain, Fusiform, Mid-Temporal}, {Ventricles}) with lower log-likelihood (log(L)=3108) than that estimated by S-EBM (log(L)=3153) but which nevertheless more closely matches the data than the serial sequence estimated by C-EBM (log(L) = 3043).

Furthermore, the C-EBM positional uncertainty can be low for groups of simultaneous events, such as T-Tau and P-Tau, demonstrating that C-EBM can be confidently incorrect regarding serial event ordering.



**Figure 6.** Positional variance diagram showing the positional uncertainty in the serial sequence estimated by C-EBM. Boxes depict the biomarkers grouped into the same stage by S-EBM. The heterogeneity within boxes indicates that C-EBM uncertainty does not infer the same information about simultaneous events as S-EBM. Furthermore, C-EBM can be confidently incorrect regarding serial orderings.

17

## 5. Conclusion

This study introduces the simultaneous event-based model. S-EBM is a generalisation of the conventional event-based model for estimating disease progression patterns that contain simultaneous events. With moderate sample sizes, S-EBM produces highly accurate sequence estimates for a range of different sequence types, including serial sequences, thereby broadening the scope of event-based modelling. By removing the requirement that the number of disease progression stages correlates linearly with the number of input biomarkers, the approach suggests a simpler explanation of AD progression, with biomarker abnormality occurring simultaneously within biomarker domains. S-EBM may provide new insights into disease evolution and more accurate subject staging, facilitating the development of therapeutic interventions targeting early disease.

## Acknowledgments

## 6. References

1. Braak, H. and Braak, E. 1991. Neuropathological stageing of Alzheimer-related changes. *Acta neuropathologica*, 82(4)

2. Crane, P.K., Carle, A., Gibbons, L.E., Insel, P., Mackin, R.S., Gross, A., Jones, R.N., Mukherjee, S., Curtis, S.M., Harvey, D. and Weiner, M. 2012. Development and assessment of a composite score for memory in the Alzheimer's Disease Neuroimaging Initiative (ADNI). *Brain imaging and behavior*, 6(4)

3. Eshaghi, A., Marinescu, R.V., Young, A.L., Firth, N.C., Prados, F., Jorge Cardoso, M., Tur, C., De Angelis, F., Cawley, N., Brownlee, W.J. and De Stefano, N. 2018. Progression of regional grey matter atrophy in multiple sclerosis. *Brain*, 141(6)

4. Fonteijn, H.M., Modat, M., Clarkson, M.J., Barnes, J., Lehmann, M., Hobbs, N.Z., Scahill, R.I., Tabrizi, S.J., Ourselin, S., Fox, N.C. and Alexander, D.C. 2012. An event-based model for disease progression and its application in familial Alzheimer's disease and Huntington's disease. *NeuroImage*, 60(3)

5. Gabel, M.C., Broad, R.J., Young, A.L., Abrahams, S., Bastin, M.E., Menke, R.A., Al-Chalabi, A., Goldstein, L.H., Tsermentseli, S., Alexander, D.C. and Turner, M.R. 2020. Evolution of white matter damage in amyotrophic lateral sclerosis. *Annals of clinical and translational neurology*, 7(5)

6. Jack Jr, C.R., Bernstein, M.A., Fox, N.C., Thompson, P., Alexander, G., Harvey, D., Borowski, B., Britson, P.J., L. Whitwell, J., Ward, C. and Dale, A.M. 2008. The Alzheimer's disease neuroimaging initiative (ADNI): MRI methods. *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine*, 27(4)

7. Jack Jr, C.R., Knopman, D.S., Jagust, W.J., Shaw, L.M., Aisen, P.S., Weiner, M.W., Petersen, R.C. and Trojanowski, J.Q. 2010. Hypothetical model of dynamic biomarkers of the Alzheimer's pathological cascade. *The Lancet Neurology*, 9(1)), 119-128

8. Marinescu, R.V., Oxtoby, N.P., Young, A.L., Bron, E.E., Toga, A.W., Weiner, M.W., Barkhof, F., Fox, N.C., Klein, S. and Alexander, D.C. 2018. Tadpole challenge: Prediction of longitudinal evolution in Alzheimer's disease. arXiv preprint arXiv:1805.03909.

9. Mueller, S.G., Weiner, M.W., Thal, L.J., Petersen, R.C., Jack, C.R., Jagust, W., Trojanowski, J.Q., Toga, A.W. and Beckett, L. 2005. Ways toward an early diagnosis in Alzheimer's disease: the Alzheimer's Disease Neuroimaging Initiative (ADNI). *Alzheimer's & Dementia*, 1(1)

10. Reuter, M., Schmansky, N.J., Rosas, H.D. and Fischl, B. 2012. Within-subject template estimation for unbiased longitudinal image analysis. *Neuroimage*, 61(4)

11. Shaw, L.M., Vanderstichele, H., Knapik-Czajka, M., Clark, C.M., Aisen, P.S., Petersen, R.C., Blennow, K., Soares, H., Simon, A., Lewczuk, P. and Dean, R. 2009. Cerebrospinal fluid biomarker signature in Alzheimer's disease neuroimaging initiative subjects. *Annals of neurology*, *65*(4)

12. Wijeratne, P.A., Young, A.L., Oxtoby, N.P., Marinescu, R.V., Firth, N.C., Johnson, E.B., Mohan, A., Sampaio, C., Scahill, R.I., Tabrizi, S.J. and Alexander, D.C. 2018. An image-based model of brain volume biomarker changes in Huntington's disease. *Annals of clinical and translational neurology*, 5(5)

13. Young, A.L., Oxtoby, N.P., Daga, P., Cash, D.M., Fox, N.C., Ourselin, S., Schott, J.M. and Alexander, D.C., 2014. A data-driven model of biomarker changes in sporadic Alzheimer's disease. *Brain*, 137(9)