

Using Historical Data to Facilitate Clinical Prevention Trials in Alzheimer Disease?

Manfred Berres (✉ berres@hs-koblenz.de)

University of Applied Sciences Koblenz: Hochschule Koblenz <https://orcid.org/0000-0003-0409-8554>

Andreas U. Monsch

University Veterinary Medicine FELIX PLATTER: Felix Platter Spital

René Spiegel

University Veterinary Medicine FELIX PLATTER: Felix Platter Spital

Research Article

Keywords: Historical controls, MCI criteria, clinical trial, cohort study, convenience sample, meta-analysis

Posted Date: February 23rd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-236517/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Using Historical Data to Facilitate Clinical Prevention Trials in Alzheimer Disease?

An Analysis of Longitudinal MCI (Mild Cognitive Impairment) Data Sets

Manfred Berres^{1,2}, Andreas U. Monsch², René Spiegel²

¹University of Applied Sciences Koblenz, Koblenz, Germany ²University Department of
Geriatric Medicine FELIX PLATTER, Basel, Switzerland

Email addresses:

berres@hs-koblenz.de,

andreas.monsch@felixplatter.ch

rene.spiegel@unibas.ch

Abstract

Background

The Placebo Group Simulation Approach (PGSA) aims at partially replacing Randomized Placebo-Controlled Trials (RPCTs) using data from historical control groups in order to decrease the needed number of study participants exposed to lengthy placebo treatment. PGSA algorithms were originally derived from Mild Cognitive Impairment (MCI) data of the Alzheimer's Disease Neuroimaging Initiative (ADNI) database. To produce more generalizable algorithms, we aimed to compile five different MCI databases in an heuristic manner to create a 'standard control algorithm' for use in future clinical trials.

Methods

We compared data from two North American cohort studies (n= 395 and 4,328, respectively) one international clinical drug trial (n= 831) and two convenience patient samples, one from Germany (n=726), and one from Switzerland (n=1,558).

Results

Despite differences between the five MCI samples regarding inclusion and exclusion criteria, their baseline demographic and cognitive performance data varied less than expected. However, the five

30 samples differed markedly with regard to their subsequent cognitive performance and clinical
31 development: (1) MCI patients from the drug trial did not deteriorate on verbal fluency over 3
32 years, whereas patients in the other samples did; (2) relatively few patients from the drug trial
33 progressed from MCI to dementia (about 10% after 4 years), in contrast to the other four samples
34 with progression rates over 30 percent.

35 *Conclusion*

36 Conventional MCI criteria were insufficient to allow for the creation of well-defined and
37 internationally comparable samples of MCI patients. More recently published MCI criteria are
38 unlikely to remedy this situation. The Alzheimer scientific community needs to agree on a standard
39 set of neuropsychological tests including appropriate selection criteria to make MCI a scientifically
40 more useful concept. Patient data from different sources would then be comparable and the
41 scientific merits of algorithm-based study designs such as the PGSA could be properly assessed.

42 **Keywords:** Historical controls, MCI criteria, clinical trial, cohort study, convenience sample,
43 meta-analysis

44 Introduction

45 Almost 10 years ago, our group published the PGSA (Placebo Group Simulation Approach)
46 for debate to the Alzheimer community [1]. The proposed novel study design was intended to
47 partially substitute for RPCTs (Randomized Placebo-Controlled Trials), i.e., clinical studies
48 which by definition expose some of the participants to treatment with placebo. We argued that,
49 in the case of Alzheimer's Disease (AD), clinical prevention trials with pre-clinical subjects
50 would typically last 18 months or longer – and that it was ethically problematic to put
51 individuals with a high risk of developing dementia on an *a priori* inactive long-term
52 medication. Instead of a concomitant placebo group the PGSA introduced algorithm-based
53 forecasts of trials subjects' expected own disease trajectories to account for the effects of
54 baseline differences, time in the study and the circumstances of trial participation. The original
55 PGSA algorithms were derived from the Alzheimer's Disease Neuroimaging Initiative (ADNI)
56 data [2] available at that time, i.e., from a then recent but, nonetheless, historical data set.

57 The pros and cons of using historical data in clinical trials have been discussed under a
58 number of aspects [3]. Thus, in the case of rare diseases it may be difficult to recruit sufficient
59 numbers of patients for proper control groups in addition to the treatment group [4]. In the case
60 of progressive, non-reversible and potentially fatal diseases there are ethical issues limiting the
61 use of inactive treatment, and it may also be difficult to obtain consent from patients for
62 participation if one of the treatment options looks more promising than the other. Finally, in
63 situations where no effective treatments are available and a promising new treatment is to be
64 tested, it may be difficult to convince patients to participate in a study which includes a placebo
65 arm [5]. However, the absence of certain subgroups of patients in a clinical trial will lead to
66 less representative samples and is then likely to cause bias in the results. Historical data could
67 be considered in some of these situations as a potential substitute for a concomitant control
68 group (for ALS see [6]).

69 A necessary prerequisite for using historical controls is that they are comparable to the
70 current study population. This requirement is hardly ever fulfilled. If only demographic
71 variables such as age, sex and observable health status are different, adjustments for these
72 covariates might solve the problem. If one attempts to set up a model for such adjustments,
73 several potential historical controls need to be compared [7].

74 In the specific case of mild cognitive impairment (MCI), a clinical condition between normal
75 aging and dementia (see definition in [8]), the problem of finding adequate historical data is
76 aggravated by the fact that MCI criteria have shifted over time and that there is no rigorous and
77 generally accepted definition of the condition [9–11]. As a typical consequence thereof,
78 different clinical drug trials with MCI subjects in the past have applied different inclusion and
79 exclusion criteria [12]. In this paper, we investigate whether information from large historical
80 MCI studies can be summarized in an heuristic manner such that a ‘standard control group’ for
81 use in future clinical trials with this population could be created. We will compare data from
82 two cohort studies, one clinical drug trial and two convenience patient samples to investigate:

- 83 • inclusion and exclusion criteria for MCI applied in five different patient datasets,
- 84 • the selection of cognitive tests applied at study entry and at follow-up,
- 85 • the homogeneity of patients’ demographic and baseline data, and
- 86 • the homogeneity of disease progression as measured by cognitive tests and indicated
87 by the proportion of conversions to dementia.

88 While distributions of demographic and baseline data will be shown, the progression of the
89 disorder is analyzed as the ‘effect of no treatment over time’. In a clinical study this
90 corresponds to the progression observed in the control group and would be contrasted to the
91 progression observed in the treatment group. In the five data sets considered, an overall
92 judgement was provided at each patient visit by an experienced clinician as to whether the

93 patient had progressed (‘converted’) from MCI to dementia. Conversion rates and hazard ratios
94 will be compared between studies.

95 In all the studies considered in this analysis, and in particular in the clinical drug trial [13],
96 some patients were treated with anti-Alzheimer medication such as cholinesterase inhibitors or
97 memantine; however, we will consider all subjects as ‘untreated’, because none of these
98 substances were shown to have significant and maintained effects on the progression of the
99 disease.

100 Material and Methods

101 *Datasets used*

102 We analyzed individual patient data from:

- 103 • the Alzheimer’s Disease Neuroimaging Initiative (ADNI;
104 <http://www.loni.ucla.edu/ADNI>)
- 105 • the National Alzheimer’s Coordinating Centers (NACC;
106 <https://www.alz.washington.edu>)
- 107 • the InDDEx clinical trial [13]
- 108 • the German Dementia Competence Network (CNG; [14])
- 109 • the Basel University Memory Clinic (BS-MC)

110 The ADNI and the NACC samples included only patients who were between 54 and 90 years
111 old at entry. BS-MC included only patients with at least 7 years of education. These restrictions
112 were applied to all 5 datasets in our analyses.

113
114 The **ADNI** (Alzheimer’s Disease Neuroimaging Initiative) study aims at investigating the
115 prognostic value of biomarkers, in particular of MRI and PET images, to describe the

116 progression of Alzheimer's disease from its preclinical to its symptomatic stages. It is led by
117 the principal investigator [2] and representatives of the ADNI sites, the NIH (National
118 Institutes of Health), the FDA (Food and Drug Administration) and contributing companies
119 from the health industry. ADNI procedures follow a detailed protocol. Cognitive performance
120 of the participants was assessed with the Alzheimer's Disease Assessment Scale – cognitive
121 subscale (ADAScog; 11 items and modified version with 13 items) [15], the MMSE [16], a
122 number of neuropsychological tests, and the Functional Assessment Questionnaire [17]. ADNI
123 started in 2003 and by the time of our last data download on January 6, 2012 [18], the data set
124 contained 395 patients diagnosed with MCI at study entry. Two subjects were excluded from
125 our analyses because they had less than 7 years of education.

126 The **NACC** (National Alzheimer's Coordinating Center) project was initiated by the National
127 Institute on Aging /NIH. It developed a large database of standardized clinical and
128 neuropathological research data collected from 29 Alzheimer's Disease Centers in the US.
129 Eight of the nine neuropsychological tests used in ADNI are also part of the NACC data base
130 [19]. We received data from the freeze of March 19, 2014. We eliminated the data from those
131 participants that were also included in the ADNI project to avoid patients from being
132 considered twice in our analysis. We selected MCI patients with memory impairment, with or
133 without impairment in other domains, who were between 55 and 90 years old and had at least 7
134 years of education. This left 4,328 MCI subjects for our analysis.

135 The **InDDEx** (Investigation of Delay of Diagnosis of AD with Exelon®) study [13] was a
136 clinical trial sponsored by Novartis Pharma, assessing the effect of the cholinesterase inhibitor
137 rivastigmine on disease progression in patients with MCI. This placebo-controlled study did not
138 show evidence of an effect of rivastigmine on either the rate of progression to dementia or the
139 standardized Z score for a cognitive test battery. We therefore considered all patients as
140 untreated and included them in our analysis. This conforms with the other four patient samples
141 where dementia-related medication was also permitted. The study applied the ADAScog, a

142 neuropsychological battery that had only verbal fluency (animals) and the Boston Naming Test
143 [20] in common with the battery used in ADNI and a different functional assessment scale. The
144 InDDEx study enrolled 1,018 patients randomly assigned to rivastigmine (n=508) and to
145 placebo (n=510). After exclusions due to missing screening data, missing cognitive data, or age
146 or education outside the admissible range, 861 subjects could be included in the present
147 analysis.

148 The **CNG** (Competence Network Germany) study [14, 21] is a longitudinal multicenter
149 cohort study of 14 memory clinics in Germany. It applies the ADAScog (12 subtests), six tests
150 of the Consortium to Establish a Registry for Alzheimer’s Disease – Neuropsychological
151 Assessment Battery (CERAD-NAB) [22] and other tests. We received data of 787 patients with
152 a diagnosis of MCI. After exclusion due to age and education restrictions 726 were left for our
153 analysis.

154 The **BS-MC** (Basel Memory Clinic) sample comprises data of patients referred by practicing
155 physicians to the memory clinic of the University Hospital Basel, Switzerland, for diagnosis
156 and treatment recommendations. Neuropsychological tests include the CERAD-NAB plus
157 Phonemic Fluency (S-words) and Trail Making Tests A and B [23] and Digit Span Forward
158 and Backward. Data of 2,135 patients with MCI at baseline were downloaded on September 9,
159 2016. After application of age and education inclusion criteria data from 1,558 patients were
160 left for the analysis.

161

162 *Statistical Analysis*

163 Demographics and baseline scores of frequently used cognitive tests are summarized in tables and
164 partly in boxplots. To investigate the homogeneity of progression across the five patient samples,
165 we performed meta-analyses for the changes from baseline of cognitive test scores. Confidence
166 intervals in forest plots will show whether there are distinct differences between studies. Measures

167 of heterogeneity confirm these results. Rates of conversion from MCI to dementia will be shown in
168 Kaplan-Meier curves, broken down by study, and hazard ratios for age, sex and education will be
169 compared in proportional hazards models.

170 Results

171 *Definition of MCI*

172 The different inclusion criteria for MCI applied in the five studies are summarized in Table 1.
173 ADNI and CNG used the MMSE, although with different inclusion criteria: The lower limit for
174 the MCI was 24/30 in ADNI, but 20/30 in CNG. These two studies requested a Clinical
175 Dementia Rating (CDR) [24] score of 0.5. Cognitive complaints or symptoms were requested
176 in ADNI, NACC and CNG. ADNI used thresholds in Logical Memory II dependent on years of
177 education. InDDEx requested a delayed recall score in the NYU delayed paragraph recall [25]
178 of less than 9, CNG and BS-MC requested at least one cognitive domain below -1 SD (CNG)
179 or -1.28 SD (BS-MC). Patients with major depression were excluded in ADNI, InDDEx and
180 BS-MC. More details on eligibility, inclusion and exclusion criteria are listed in Table 1.

181 Table 1: Inclusion and exclusion criteria for the diagnosis of MCI used in five studies.

	Eligibility			Inclusion		Exclusion	
	Age	MMSE	Other, Medication	Functional Impairment	Cognitive impairment	Depression	Other/Vascular
ADNI	55-90	24-30	Stable medication, AChEIs, memantine admitted, 6 grades education or work history	No functional impairment, but many with high FAQ scores. CDR=0.5; memory ≥ 0.5	Memory complaint LogMem II, dependent on education	Geriatric Depression Scale ≥ 6	Hachinski Ischemic Score IS >5
NACC	--	--	Similar to ADNI	Essentially normal daily functions	Cognitive complaint, cognitive decline (clinician's diagnosis)	Not specified	Not specified
InDDEx	55-85	--	No AChEI in previous 2 weeks, no rivastigmine in previous 4 weeks	Cognitive symptoms (not specified); CDR=0.5	NYU delayed paragraph recall <9	HDRS >12 , HDRS item 1 > 1 , DSM-IV major depression	AD criteria from DSM-IV or NINCDS-ADRDA mod. Hachinski Ischemic Score >4
CNG	≥ 50	≥ 20	A broader definition of MCI was used	Complaint of cognitive deficit in daily living; minor changes were tolerated: B-ADL < 4	Decline of cog. abilities (>1 SD) in at least one neuropsychological domain	Not specified	Not specified
BS-MC	N/A	N/A	Consecutively referred patients from GPs	Essentially Winblad et al. (2004)[26] criteria; no significant functional decline	Impairment (≤ -1.28 SD; age-, education- and gender-adjusted) in \geq one cognitive domain	Probable cause for MCI other than early AD, based on comprehensive medical exam and neuroimaging results	Not specified

182

183 *Cognitive tests*

184 A selection of tests, most of them applied in at least two studies, is listed in Table 2. Each study
185 applied a different set of tests. The Mini Mental Status Examination (MMSE) and the Verbal
186 Fluency Test (animals) were the only instruments used in all studies. The ADAScog [15] was
187 applied, although with different modifications, in ADNI, InDDEx and CNG. Eight of nine tests of a
188 neuropsychological battery used in ADNI were applied in NACC as well, but the Auditory Verbal
189 Learning Test was omitted. Moreover, several procedural details of the Logical Memory delayed
190 recall from the Wechsler Memory Scale (WMS) differed significantly in ADNI and NACC (details

191 in [18]), The Boston Naming Test [20] was performed in ADNI and NACC with 30 items, but with
192 only 15 items in CNG and BS-MC. The Digit span forward and backward and the Trail Making
193 Test A and B [27] were applied in all studies except in InDDEx; however, the Trail Making Tests
194 were conducted with different time limits (e.g., for TMTA: 150 seconds in NACC and CNG, 180
195 seconds in BS-MC). The Clock Drawing Test was applied in all studies except in NACC, but
196 different scoring methods were used. The Clinical Dementia Rating scale was applied in all studies
197 but BS-MC. Three different versions of functional assessment were used. The CERAD battery was
198 only used in CNG and BS-MC. A few other tests were applied in only one study.

199 Table 2: List of selected cognitive tests applied in five studies

Test	ADNI	NACC	InDDEx	CNG	BS-MC
ADAScog 11 and modified	11 & mod.		11 & mod	11	
Logical Memory II	x	x		x	
Digit Span Forward	x	x		x	x
Digit Span Backward	x	x		x	x
Category Fluency, Animals	x	x	x	x	x
Category Fluency, Vegetables	x	x			
Trail Making Test B	x	x		x	x
Boston Naming Test	x	x	x*	x	x
Auditory Verbal Learning Test	x				
Digit Symbol	x	x			
Trail Making Test A	x	x		x	x
Clock Drawing Test	x		x	x	x
Functional Assessment	x	x	ADCS-ADL	Bayer-ADL	
Logical Memory I	x			x	
American National Adult Reading Test	x				
Clinical Dementia Rating	x	x	x	x	
Mini Mental Status Examination	x	x	x	x	x
Phonemic fluency, S-words					x
CERAD Wordlist + intrusion errors / savings				x	x
CERAD constructional praxis				x	x
Neuropsychiatric Inventory (NPI)	x	x	x	x	
Digit cancellation task			x		

200 * only 85 values

201

202

203 *Demographics (Table 3)*

204 Patients in ADNI and NACC were oldest (74.2 ± 7.5 and 74.4 ± 7.9), those in InDDEx and BS-MC
 205 were younger (70.1 ± 7.6 and 69.7 ± 9.1), CNG (68.0 ± 7.9) comprised the youngest sample. Duration
 206 of education was highest in ADNI (15.7 ± 3.0 years) and NACC (15.2 ± 3.0 years) and lowest in
 207 InDDEx (11.8 ± 3.5). CNG (12.3 ± 2.8) and BS-MC (12.0 ± 3.1) were in between. The proportion of
 208 females was lowest in ADNI (35.7%), highest in NACC (48.8%) and InDDEx (49.0%) and
 209 intermediate in CNG (45.7%) and BS-MC (45.8%). The proportion of patients with two ApoE4
 210 alleles was highest in ADNI (11.9%), intermediate in NACC (8.7%), InDDEx (9.6%) and BS-MC
 211 (9.4%) and lowest in CNG (5.4%).

212 Table 3: Descriptive statistics of demographics and baseline scores MMSE, Verbal Fluency
 213 (animals) and ADAScog (11 subtests).

	ADNI n = 395	NACC n = 4,328	InDDEx n = 831	CNG n = 726	BS-MC n = 1,558
Age ($\bar{x} \pm SD$)	74.2±7.5	74.4±7.9	70.1±7.6	68.0±7.9	69.7±9.1
Education ($\bar{x} \pm SD$)	15.7±3.0	15.2±3.0	11.8±3.5	12.3±2.8	12.0±3.1
Female n (%)	141 (35.7)	2113 (48.8)	422 (49.05)	332 (45.7)	713 (45.8)
ApoE4*	n _{ApoE} =395	n _{ApoE} =2,523	n _{ApoE} =396	n _{ApoE} =577	n _{ApoE} =53
1 allele no (%)	165 (41.8)	951 (37.7)	131 (32.3)	200 (34.7)	22 (41.5)
2 alleles no (%)	47 (11.9)	219 (8.7)	39 (9.6)	31 (5.4)	5 (9.4)
MMSE ($\bar{x} \pm SD$)	27.0±1.8	27.0±2.4	27.2±2.5	27.1±2.1	27.4±2.3
min-max	23-30	2-30	16-30	17-30	14-30
Verb. Fl. ($\bar{x} \pm SD$)	15.9±4.9	16.0±5.0	17.5±5.9	17.5±5.5	17.4±5.5
min-max	5-30	0-35	2-38	3-32	3-38
ADAScog ($\bar{x} \pm SD$)	11.5±4.4	-	10.0±4.7	11.7±5.1	-
min-max	2-27.7		1-27	0-35	

214 * ApoE was not determined in all patients, number of evaluations is given as n_{ApoE}, percent of ApoE
 215 evaluations shown in parentheses.

216
217 *Baseline data*

218 In each study the quartile range of MMSE is from 26 to 29 (Figure 1), but only ADNI had no
219 outliers, because the minimum score of 24 was an inclusion criterion. Verbal fluency is 1.4 to 1.6
220 points higher in InDDEx, CNG and BS-MC compared to ADNI and NACC (Table 3, Figure 2).
221 InDDEx patients performed better on the ADAScog (10.0 ± 4.7) than CNG patients (11.7 ± 5.1), but
222 ADNI patients (11.5 ± 4.4) performed similar to CNG patients. In the Boston Naming Test [20]
223 patients of ADNI and NACC achieved on average about 25 of 30 items, while patients of CNG and
224 BS-MC achieved about 13.4 of 15 items. Results on the Trail Making test cannot be compared
225 because of different time limits used.

226 *Place Figure 1 about here*

227 *Place Figure 2 about here*

228 *Progression of cognitive scores*

229 Verbal Fluency (animals) is – next to the MMSE – the only test score available in all five studies.
230 From baseline to one year, InDDEx patients improved on average by 0.4 words, CNG patients
231 remained stable, but ADNI, NACC and BS-MC patients worsened by 0.6, 0.7 and 0.8 words,
232 respectively (Figure 3). Confidence intervals for InDDEx and the latter three studies are disjoint.
233 The differences between studies increased for the two- and three-year follow-up: InDDEx subjects
234 still improved by 0.5 words, the latter three worsened by up to 2.1 (NACC) and 3.5 (BS-MC) words
235 after 3 years (Figure 3).

236 *Place Figure 3 about here*

237 ADAScog worsened considerably in ADNI, slightly in CNG, but improved slightly in InDDEx.
238 Boston Naming Test scores show similar decline in: ADNI, NACC and BS-MC (in relation to the
239 number of items), while the decline in CNG is very small. The CERAD-word list (delayed recall)
240 worsened considerably in BS-MC, but remained unchanged in CNG.

241 *Conversion from MCI to dementia*

242 Time to progression (conversion) from MCI to dementia was usually ascertained at scheduled
243 visits. Sometimes, however, patients were examined in an extra visit and then classified as
244 dementia. Kaplan-Meier curves for the time to conversion from MCI to dementia (mostly AD) are
245 shown in Figure 4. They confirm that the InDDEx patients were in a particularly stable state.
246 Continuation with an MCI diagnosis was considerably less frequent in the CNG sample, while
247 ADNI and NACC patients were the fastest to progress. Conversion rates after 3 years are between
248 5.9% (InDDEx) and 46.4% (NACC), both with a standard error of 1.1%

249 *Place Figure 4 about here*

250 Cox regression models with covariates years of education, age and sex were estimated for each
251 study. Hazard ratios and 95% confidence intervals are shown in Figure 5. The hazard ratio for
252 education was 0.71 (for 4 years) in CNG and close to 1 in the other studies. The hazard ratio for age
253 was close to 1 in ADNI and distinctly positive in the other studies. The hazard ratio for females was
254 1.56 in BS-MC and close to 1 in the other studies.

255 *Place Figure 5 about here*

256 Discussion and Conclusion

257 The PGSA (Placebo Group Simulation Approach) was submitted ‘*for debate*’ to the Alzheimer
258 community [1]. The novel study design was intended primarily to resolve an ethical problem of
259 prevention trials involving aged subjects at risk of shifting into dementia: the long-term use of
260 placebo typical of RCTs (Randomized Controlled Trials). Accordingly, instead of using a
261 concomitant placebo group for comparison with a hopefully effective novel treatment, the PGSA
262 applies mathematical algorithms to forecast the expected outcomes of pre-symptomatic AD patients
263 from their baseline data and to compare those with the outcomes on experimental treatments. The
264 PGSA was deemed to “have an advantage over the use of historical controls in futility designs in
265 that it is based on the patient’s own observed clinical features.” [28].

266 Our published algorithms were derived from the ADNI database (<http://www.loni.ucla.edu/ADNI>)
267 [29] that contained anamnestic, biological, neuroimaging and neuropsychological findings from 397
268 North American patients with a diagnosis of MCI (Mild Cognitive Impairment). Our analyses
269 highlighted the strong impact of neuropsychological performance data recorded at baseline, in
270 addition to information from subjects' history such as age, sex and education, on MCI disease
271 trajectories in the following years. A first attempt at validation of the PGSA algorithms using data
272 from an independent MCI database (NACC; <https://www.alz.washington.edu>) confirmed the
273 importance of neuropsychological performance data recorded at baseline to forecast cognitive
274 decline in MCI [18]. However, we also noted that there was some slight over-estimation of
275 cognitive decline when the ADNI-based PGSA algorithms were applied to the NACC MCI dataset
276 for a follow-up of more than two years. This observation led to the question as to whether the
277 published PGSA algorithms could be confidently applied to other longitudinal MCI data.

278 The current analysis comprised three MCI databases in addition to the ones from ADNI and
279 NACC. One of these three originated from an RCT with a cholinesterase inhibitor sponsored by a
280 drug company [13] and two from clinical case series: one from a network of memory clinics in
281 Germany [21], the other one from a single memory clinic in Basel, Switzerland [30]. The five
282 databases contained information on 395 (ADNI) up to 4328 (NACC) individuals. Although all
283 patients had a diagnosis of MCI, inspection of Table 1 shows that inclusion and exclusion as well as
284 other eligibility criteria varied considerably between the five samples. Wide differences also existed
285 between the five databases with regard to the neuropsychological scales and instruments used at
286 baseline and at follow-up to document the changes in patients' cognitive performance (Table 2).
287 Only the MMSE and the Verbal Fluency test with "animals" were applied throughout.

288 Despite the differences in inclusion and exclusion criteria, the demographic and even more so the
289 baseline cognitive performance data of the five patient samples were not as variable as one might
290 expect. While patients' mean ages varied between 68.0 (CNG) and 74.4 (NACC) years and the
291 ADNI sample contained a lower percentage of female participants than the other four, mean MMSE

292 scores at baseline varied only minimally (Table 3, Figure 1), and the same was seen with regard to
293 the ADAS scores in the three studies where this scale was used. Interestingly, the participants of
294 InDDEx, CNG and BS-MC had slightly better performance on the Verbal Fluency test than those of
295 ADNI and NACC (Table 3, Figure 2), although the latter had benefited from more years of
296 education on the average.

297 In contrast to the similarity of their cognitive performance data noted at baseline, the five samples
298 differed markedly with regard to their subsequent cognitive performance and clinical development.
299 As seen in Figure 3, scores on the Verbal Fluency test behaved differently in the InDDEx than in
300 the other four groups: with increasing study duration performance deteriorated continuously in the
301 ADNI, NACC, CNG and BS-MC samples, but were slightly improved from baseline after one year
302 and then stayed stable in the InDDEx group. An even greater disparity is seen with regard to the
303 number of conversions from MCI to dementia (Figure 4): Whereas some 90 percent of the InDDEx
304 patients did not progress to dementia within the 4 years of observation, the respective percentages
305 were around 70 for CNG, around 60 for ADNI and BS-MC and somewhat above 50 for NACC.
306 What could be an explanation of these big differences?

307 First, one should note that the difference in conversion rates between the ADNI and the NACC
308 participants is small. This is not a surprise given that both MCI patient samples were collected in
309 North America, and that the ADNI sample partly constituted a selection from the large NACC data
310 collection. Thus, it is likely that both the MCI inclusion/exclusion and the conversion criteria
311 applied to the ADNI and the NACC data were similar. Progression rates for the BS-MC sample
312 were also close to the ones seen in ADNI and NACC, suggesting that the inclusion/exclusion
313 criteria for MCI and the criteria to diagnose conversion to dementia were applied similarly in the
314 Basel and in the North American centers. More difficult to understand are the lower conversion
315 rates in CNG and, particularly, in InDDEx. The latter dataset differed from the other four in that
316 InDDEx was not a case series from one or more memory clinics but a clinical drug trial sponsored
317 by a pharmaceutical company and carried out in 12 different countries in Europe, South Africa,

318 South America and in the USA. While one cannot exclude that this geographic variety (or perhaps
319 some investigators' desire to include as many patients as possible) compromised proper selection of
320 MCI patients, it is of note that other, although shorter, company-sponsored drug studies carried out
321 in the same decade as InDDEx also differed markedly with regard to conversion rates. Thus,
322 Petersen [8] reported annual conversion rates of 16 percent in a 3-year study with donepezil and
323 Vitamin E , and relatively high percentages of conversions were also seen in two separate 2-year
324 studies with galantamine [31]. In contrast, an RCT with a selective COX-2 inhibitor noted rather
325 low annual conversion rates: 6.4% on rofecoxib, 4.5% on placebo [32], - although these conversion
326 rates were still higher than the one reported in [13].

327 Whatever the reason for these varying conversion rates were, it is obvious that the MCI criteria
328 available some 20 years ago were not sufficiently precise to allow selection of clinically
329 homogeneous and internationally comparable MCI patient groups. The problem was noted early on
330 [9, 33] and efforts were made to ameliorate the situation (e.g. [26]). Nevertheless, ambiguity
331 remained: Although an international expert group [10] stressed the importance of determining
332 whether there is objective evidence of cognitive decline and recommended cognitive testing for
333 quantitatively assessing the degree of cognitive impairment for a diagnosis of MCI, these authors
334 also emphasized that normative ranges of neuropsychological tests (typically those listed in Table
335 2) "are guidelines and not cutoff scores." ([10] p. 272). In contrast to this position, it is our opinion
336 that the scientific community needs to decide on age- and education-adjusted cutoff scores in order
337 to make MCI a scientifically useful concept and to ensure that study results from different sources
338 will be comparable. Another critical issue is the use of specific neuropsychological tests: as noted in
339 the current analyses only two tests were part of all five studies considered, making comparison
340 between patient samples virtually impossible. Moreover, for several tests different versions and
341 scoring systems exist and are being used. This is another impediment when one tries to compare
342 results between studies. In the USA, a series of Uniform Data Sets (UDS) have been proposed in
343 [34] to improve this unsatisfactory situation. A separate issue refers to the underlying disorder of

344 MCI and subsequent dementia. While neuropsychological test performance with an internationally
345 agreed set of tests and cutoffs will allow for the determination of cognitive deficits, specific
346 biomarker requirements would also allow to describe pathophysiologically more homogenous
347 groups of patients with MCI, e.g., MCI due to AD. As of today, and as shown in our analyses, the
348 heterogeneity at all these levels is way too big to allow meaningful conclusions with regard to the
349 scientific merits of algorithm-based approaches such as the PGSA.

350 The scientific rationale of our earlier studies [1, 15] was to make use of well-defined historical
351 data in clinical treatment or prevention trials. Inclusion of historical information for comparison
352 with current treatment data has been discussed since more than 40 years ([35] and later references in
353 [7]). Suggestions range from performing a single arm study which is compared to the historical
354 control, over integrating historical control data with new controls - up to ignoring historical data
355 completely. Different options are available to integrate historical data [3]: (1) pooling with new
356 controls; (2) testing for differences between historical and new control and pool only if no
357 differences are detected; (3) down weighting the historical data by power priors dependent on the
358 discrepancy between observed and historical control data; (4) choose a prior distribution for the
359 means of the historical and the new control groups and apply a Bayesian model (see [3] for details
360 on these methods); (5) perform a random effects meta analysis of the historical controls and down-
361 weight their sample size according to the between-study variation [7].

362 However, data pooling is only admissible if the historical controls are exactly equivalent to the
363 new control. This may hold for a set of pharmaceutical studies with basically the same protocol and
364 equivalent patient populations. In less narrowly defined circumstances this is rarely the case, due to
365 different populations and more or less different inclusion and exclusion criteria. Down-weighting of
366 historical data takes heterogeneity into account. It decreases the weight of the historical data from
367 the total number of patients to a smaller number which is called the *prior effective sample size* [7].
368 For the five studies considered in the current analysis, the 2,144 patients for whom data for change
369 of verbal fluency after three years were available, would be down weighted to 10 patients,

370 according to formulas [7]. However, all methods of integrating historical controls are only valid
371 under the assumption of exchangeability, i.e., if no systematic differences exist between the control
372 groups [7, 36]. In view of the distinctly apart confidence intervals in Figure 3, exchangeability
373 cannot be assumed for the studies considered here. This makes all attempts to integrate these data in
374 new studies futile. Insofar, our current failure of forecasting algorithms is consistent with theoretical
375 consideration on using prior information.

376 *Limitations*

377 Drop-outs are a common problem in long-term clinical studies. In the five studies considered, the
378 rate of drop-outs after 3 years was quite different: 70% in NACC, 37% in ADNI, 25% in InDDEx,
379 77% in CNG and 98% in BS-MC. While studies with a stricter visit regimen (in our case InDDEx
380 and ADNI) have lower drop-out, convenience samples such as CNG and BS-MC tend to have very
381 high drop-out rates. This factor may cause bias and could make studies less comparable. For
382 example, one might assume that convenience samples contain a larger number of frail patients than
383 samples in controlled studies, that frail patients drop out with higher probability and that, as a
384 consequence, the results of cognitive tests would worsen less in convenience samples. However, our
385 results do not support this hypothesis: Patients of the InDDEx study (with the lowest dropout rate)
386 improved their cognitive results after 2 and 3 years, whereas patients in BS-MC showed the most
387 pronounced decrease (cf. Figure 3). CNG and NACC, two studies that collected data from
388 Alzheimer's coordinating centers and/or memory clinics, also showed distinctly different changes
389 (Figure 3). Conversion rates in the InDDEx sample were very low, and the conversion rate of the
390 CNG sample was between InDDEx and the other studies (cf. Figure 4). This conforms to the
391 cognitive test results, but, again, does not support the hypothesis concerning bias due to drop-out.

392 Conversions were in most studies ascertained at more or less strictly scheduled visits. We
393 nevertheless applied Kaplan-Maier and Cox Regression analysis assuming continuous time. Figure
394 4 shows the pattern of event times for strict visit schedules of ADNI and InDDEx and less clearly

395 for NACC and CNG. The true curves would interpolate between the top right bends of the curves in
396 Figure 5, but this would not change the interpretation of the conversion rates.

397

399 **References**

- 400 1. Spiegel R, Berres M, Miserez AR, Monsch AU. For debate: substituting placebo controls in
401 long-term Alzheimer's prevention trials. *Alz Res & Ther.* 2011;3:9–20.
- 402 2. Weiner MW, Veitch DP, Aisen PS, Beckett LA, Cairns NJ, Green RC, Harvey D, Jack CR,
403 Jagust W, Liu Enchi, Morris JC, Petersen RC, Saykin AJ, Schmidt ME, Shaw L, Siuciak JA,
404 Soares H, Toga AW, Trojanowski JQ. The Alzheimer's Disease Neuroimaging Initiative: A
405 review of papers published since its inception. *Alzheimer's & Dementia.* 2011;7:1–67.
- 406 3. Viele K, Berry S, Neuenschwander B, Amzal B, Chen F, Enas N, Hobbs B, Ibrahim JG,
407 Kinnersley N, Lindborg S, Micallef S, Roychoudhury S, and Thompson L. Use of historical
408 control data for assessing treatment effects in clinical trials. *Pharmaceutical Statistics.*
409 2014;13:41–54.
- 410 4. Miller RG, Moore DH, Forshew DA, Katz JS, Barohn RJ et al. Phase II screening trial of
411 lithium carbonate in amyotrophic lateral sclerosis. *Neurology.* 2011;77:973–9.
- 412 5. Grill JD & Karlawish J. Addressing the challenges to successful recruitment and retention in
413 Alzheimer's disease clinical trials. *Alzheimer's Research & Therapy.* 2010;2:34–44.
- 414 6. Cudkovicz M.E., Katz J., Moore D.H. O'Neill G., Glass J. D. et al. Toward more efficient
415 clinical trials for amyotrophic lateral sclerosis. *Amyotrophic Lateral Sclerosis.* 2010;11:259–65.
- 416 7. Neuenschwander B, Capkun-Niggli G, Branson M, and Spiegelhalter D. Summarizing
417 historical information on controls in clinical trials. *Clinical Trials.* 2010;7:5–18.
- 418 8. Petersen RC, Thomas RG, Grundman M, Bennett D, Doody R et al. Vitamin E and donepezil for
419 the treatment of mild cognitive impairment. *New Engl J Med.* 2005;352:2379–88.
- 420 9. Visser PJ, Scheltens P, Verhey FRJ. Do MCI criteria in drug trials accurately identify subjects
421 with predementia Alzheimer's disease? *J Neurol Neurosurg Psychiatry.* 76;2005:1348–54.
- 422 10. Albert MS, DeKosky ST, Dickson D, Dubois B, Feldman HH. The diagnosis of mild cognitive
423 impairment due to Alzheimer's disease: Recommendations from the National Institute on
424 Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimers disease.
425 *Alzheimer's & Dementia.* 2011;7:270–9.
- 426 11. Morris JC. Revised criteria for mild cognitive impairment may compromise the diagnosis of
427 Alzheimer disease dementia. *Arch Neurol.* 2012;69:700–8. doi:10.1001/archneurol.2011.3152.
- 428 12. Petersen R.C., Thomas R.G., Aisen P.S., Mohs R.C., Carrillo M.C. et al. Randomized
429 controlled trials in mild cognitive impairment. Sources of variability. *Neurology.*
430 2017;88:1751–8.
- 431 13. Feldman HH, Ferris S, Winblad B, Sfikas N, Mancione L et al. Effect of rivastigmine on delay
432 to diagnosis of Alzheimer's disease from mild cognitive impairment: the InDDEX study. *Lancet*
433 *Neurology.* 2007;6:501–12.
- 434 14. Kornhuber J, Schmidtke K, Frölich L, Perneckzy R, Wolf S et al. Early and Differential
435 Diagnosis of Dementia and Mild Cognitive Impairment. *Dementia Geriatric Cognitive*
436 *Disorder.* 2009;27:404–17.
- 437 15. Rosen WG, Mohs RC, Davis KL. A new rating scale for Alzheimer's disease. *Am J Psychiatry.*
438 1984;141:1356–64.
- 439 16. Folstein MF, Folstein SE, McHugh PR. "Mini-mental state". A practical method for grading the
440 cognitive state of patients for the clinician. *J Psychiatr Res.* 1975;12:189–98. doi:10.1016/0022-
441 3956(75)90026-6.
- 442 17. Pfeffer RI, Kurosaki TT, Harrah CH Jr, Chance JM, Filos S. Measurement of functional
443 activities in older adults in the community. *J Gerontol.* 1982;37:323–9.
- 444 18. Berres M, Kukull WA, Miserez RA, Monsch AU, Monsell SE, Spiegel R for the Alzheimer's
445 Disease Neuroimaging Initiative. A Novel Study Paradigm for Long-term Prevention Trials in
446 Alzheimer Disease: The Placebo Group Simulation Approach (PGSA). Application to MCI data
447 from the NACC database. *Journal of Prevention of Alzheimer's Disease.* 2014;1:99–109.

- 448 19. Weintraub S, Salmon DS, Mercaldo N, Ferris S, Graff-Radford NR, Chui H, Cummings J,
449 DeCarli Ch, Foster NL, Galasko D, Peskind E, Dietrich W, Beekly DL, Kukull WA, Morris JC.
450 The Alzheimer's Disease Centers' Uniform Data Set (UDS) The Neuropsychologic Test
451 Battery. *Alzheimer Dis Assoc Disord*. 2009;23:91–101.
- 452 20. Kaplan E, Goodglass H, Weintraub S. *The Boston naming test*. Philadelphia: Lea & Febiger;
453 1983.
- 454 21. Wolfsgruber S, Wagner M, Schmidtke K, Frölich L, Kurz A, Schulz S, et al. Memory concerns,
455 memory performance and risk of dementia in patients with mild cognitive impairment. *PLoS*
456 *ONE* 2014. doi:10.1371/journal.pone.0100812.
- 457 22. Morris JC, Heyman A, Mohs RC, Hughes JP, van Belle G, Fillenbaum G, Mellits ED, Clark C.
458 . The Consortium to Establish a Registry for Alzheimer's Disease (CERAD): I. Clinical and
459 neuropsychological assessment of Alzheimer's disease. *Neurology*. 1989;39:1159–65.
460 doi:10.1212/wnl.39.9.1159.
- 461 23. Schmid NS, Ehrensperger MM, Berres M, Beck IR, Monsch AU. The Extension of the German
462 CERAD Neuropsychological Assessment Battery with Tests Assessing Subcortical, Executive
463 and Frontal Functions Improves Accuracy in Dementia Diagnosis. *Dement Geriatr Cogn Dis*
464 *Extra*. 2014;4:322–34. doi:10.1159/000357774.
- 465 24. Morris JC. . The Clinical Dementia Rating (CDR): Current version and scoring rules.
466 *Neurology*. 1993;43:2412–4.
- 467 25. Kluger A, Ferris SH, Golomb J, Mittelman MS, Reisberg B. Neuropsychological prediction of
468 decline to dementia in nondemented elderly. *J Geriatr Psychiatry Neurol*. 1999;12:168–79.
469 doi:10.1177/089198879901200402.
- 470 26. Winblad B, Palmer K, Kivipelto M, Jelic V, Fratiglioni, L. et al. Mild cognitive impairment –
471 beyond controversies, towards a consensus: report of the International Working Group on Mild
472 Cognitive Impairment. *Journal of Internal Medicine*. 2004;256:240–6.
- 473 27. War Department, Adjutant General's Office. *Army Individual Test Battery: Manual of*
474 *directions and scoring*. Washington, DC; 1944.
- 475 28. Cummings J., Gould H. & Zhong K. Advances in designs for Alzheimer's disease clinical
476 trials. *Amer J Neurodeger Dis*. 2012;1:205–16.
- 477 29. Alzheimer's Disease Neuroimaging Initiative. 2017. <http://www.loni.ucla.edu/ADNI>. Accessed
478 09-14-2020.
- 479 30. Monsch AU, KRW. Specific care program for the older adults: Memory Clinics. *European*
480 *Geriatric Medicine*. 2010;1:28–31.
- 481 31. Winblad B, Gauthier S, Scinto L, Feldman H, Wilcock GK et al. Safety and efficacy of
482 galantamine in subjects with mild cognitive impairment. *Neurology*. 2008;70:2024–35.
- 483 32. Thal LJ, Ferris, SH, Kirby L, Block GA, Lines ChR et al. A randomized double-blind, study of
484 rofecoxib in patients with mild cognitive impairment. *Neuropharmacology*. 2005;30:1204–15.
- 485 33. Visser PJ, Kester A, Jolles J & Verhey F. Ten-year risk of dementia in subjects with mild
486 cognitive impairment. *Neurology*. 2006;67:1201–7.
- 487 34. Besser L, Kukull W, Knopman DS, Chui H, Galasko D, Weintraub S, Jicha G, Carlsson C,
488 Burns J, Quinn J, Sweet RA, Rascovsky K, Teylan M, Beekly D, Thomas G, Bollenbeck M,
489 Monsell S, Mock C, Zhou XH, Thomas N, Robichaud E, Dean M, Hubbard J, Jacka M,
490 Schwabe-Fry K, Wu J, Phelps C, Morris JC, Neuropsychology Work Group, Directors, and
491 Clinical Core leaders of the National Institute on Aging-funded US Alzheimer's Disease
492 Centers. Version 3 of the National Alzheimer's Coordinating Center's Uniform Data Set.
493 *Alzheimer Dis Assoc Disord*. 2018;32:351–8. doi:10.1097/WAD.0000000000000279.
- 494 35. Pocock S. combination of randomized and historical controls in clinical trials. *J Chron Dis*.
495 1976;29:175–88.
- 496 36. Schmidli H, Gsteiger S, Roychoudhury S, O'Hagan A, Spiegelhalter D, Neuenschwander B.
497 Robust meta-analytic-predictive priors in clinical trials with historical control information.
498 *Biometrics*. 2014;70:1023–32. doi:10.1111/biom.12242.
- 499

501 **Declarations**

502 **Ethics approval**

503 Ethics approval for using the data of the University Departement of Geriatric Medicine FELIX-
504 PLATTER has been obtained from the Ethics Committee (Ethikkommission Nordwest- und
505 Zentralschweiz). For the other datasets ethics approval is a prerequisite of these studies.

506 **Consent for publication**

507 Is not applicable beyond ethics approval.

508 **Availability of data**

509 The data that support the findings of this study are available from Alzheimer's Disease
510 Neuroimaging Initiative (ADNI; <http://www.loni.ucla.edu/ADNI>), National Alzheimer's
511 Coordinating Centers (NACC; <https://www.alz.washington.edu>), Novartis AG, Basel, Competence
512 Network Germany [21] and University Departement of Geriatric Medicine FELIX-PLATTER
513 (A.U.M) but restrictions apply to the availability of these data, which were used under license for
514 the current study, and so are not publicly available. Data are however available from the authors
515 upon reasonable request and with permission of third parties mentioned before.

516 **Competing interests**

517 None of the authors have any competing interests.

518 **Funding**

519 This study was supported by grants from the Alzheimer's Association Switzerland and the
520 Alzheimer Forum Switzerland.

521 **Authors' contribution**

522 M.B. performed the data analysis and contributed Material and Methods, Results, Limitations.
523 R.S. contributed the Introduction and most of the Discussion, A.U.M. added to the Abstract and the
524 Discussion. All authors approved the final manuscript.

525 Acknowledgment

526 Data collection and sharing for this project was funded by the ADNI (National Institutes of Health
527 Grant U01 AG024904). The ADNI is funded by the National Institute on Aging, the National
528 Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the
529 following: Abbott, AstraZeneca AB, Bayer Schering Pharma AG, Bristol-Myers Squibb, Eisai
530 Global Clinical Development, Elan Corporation, Genentech, GE Healthcare, GlaxoSmithKline,
531 Innogenetics, Johnson and Johnson, Eli Lilly and Co., Medpace, Inc., Merck and Co., Inc., Novartis
532 AG, Pfizer Inc., F. Hoffman-La Roche, Schering-Plough, Synarc, Inc., as well as nonprofit partners
533 the Alzheimer's Association and the Alzheimer's Drug Discovery Foundation, with participation
534 from the US Food and Drug Administration. Private-sector contributions to the ADNI are facilitated
535 by the Foundation for the National Institutes of Health. The grantee organization is the Northern
536 California Institute for Research and Education, and the study is coordinated by the Alzheimer's
537 Disease Cooperative Study at the University of California, San Diego. ADNI data are disseminated
538 by the Laboratory for Neuro Imaging at the University of California, Los Angeles. This research
539 was also supported by NIH grants P30 AG010129, K01 AG030514, and the Dana Foundation. The
540 support provided by Michael Wagner and Steffen Wolfsgruber from the Competence Network
541 Germany (CNG), by Walter Kukull and Sarah Monsell from the National Alzheimer Competence
542 Center (NACC) and by Peter Quarg from Novartis, Basel at an earlier stage of our project is
543 gratefully acknowledged.

544

545 *Figure titles and legends*

546

547 Figure 1: Boxplot of Mini Mental Status Examination (MMSE) scores at baseline in each study.

548 Figure 2: Boxplot of Verbal Fluency scores (animals) at baseline in each study.

549 Figure 3: Forest plots for the change of Verbal Fluency (animals) from baseline to 1, 2 and 3
550 years.

551 Legend to Figure 3: Mean changes and 95% confidence intervals for each study and for the
552 overall effect in the fixed effects and the random effects model are given. τ^2 is the between-study
553 variance, I^2 measures heterogeneity (between study variance over total variance), p -value for the
554 test of heterogeneity. Graphs show study specific means and confidence intervals for each study as
555 grey squares and lines and for overall effects as diamonds. Size of squares represent precision of
556 individual treatment estimates.

557 Figure 4: Kaplan-Meier plots of the proportion converted from MCI to dementia versus time for
558 each study.

559 Figure 5: Hazard ratios of five Cox proportional hazard regression models for each study.

560 Legend to Figure 5: For education the hazard ratio for progression to dementia is shown for an
561 increase of 4 years, for age it is shown for an increase of 10 years, for gender it is for females
562 relative to males.

Figures

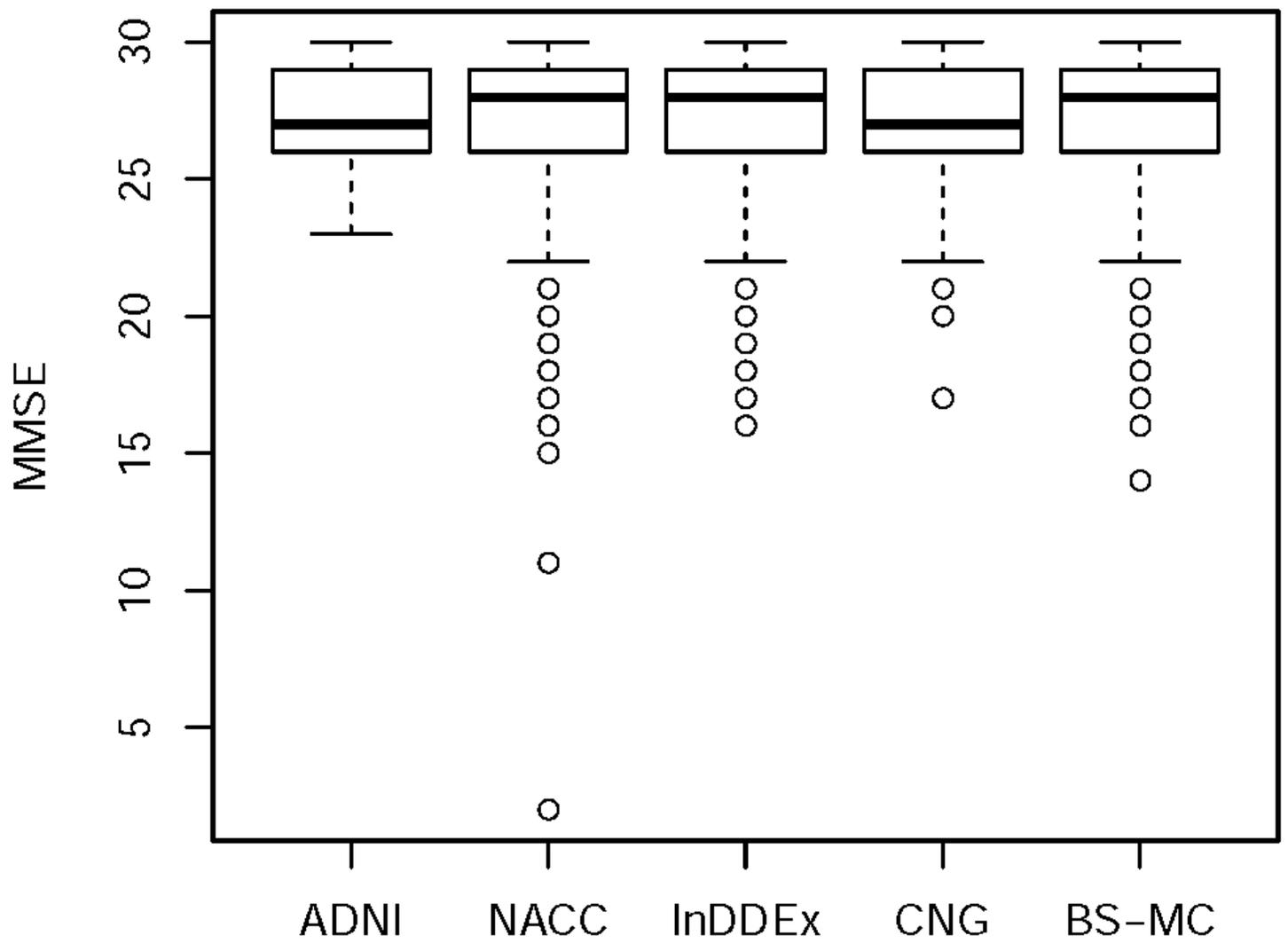


Figure 1

Boxplot of Mini Mental Status Examination (MMSE) scores at baseline in each study.

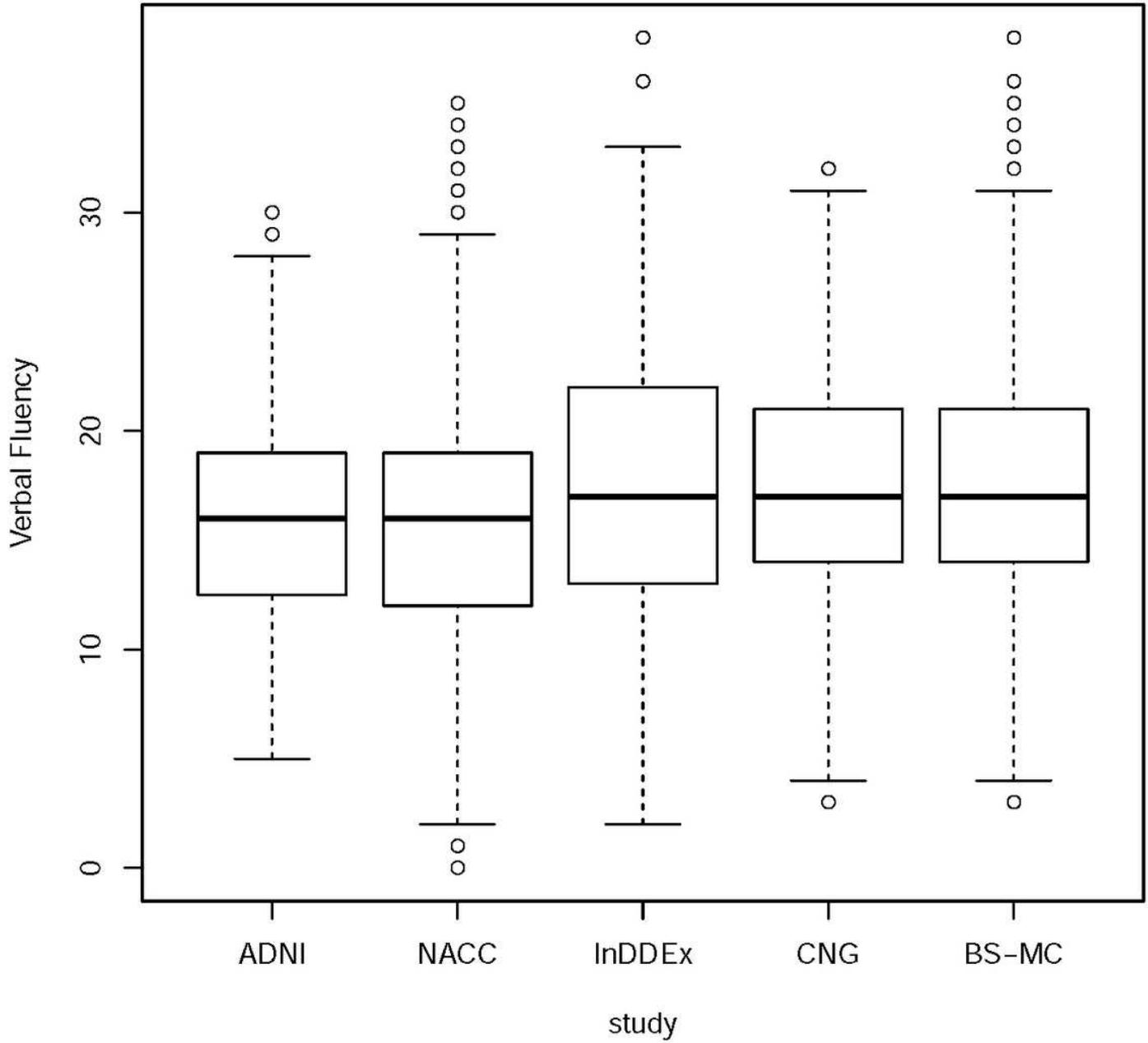


Figure 2

Boxplot of Verbal Fluency scores (animals) at baseline in each study.

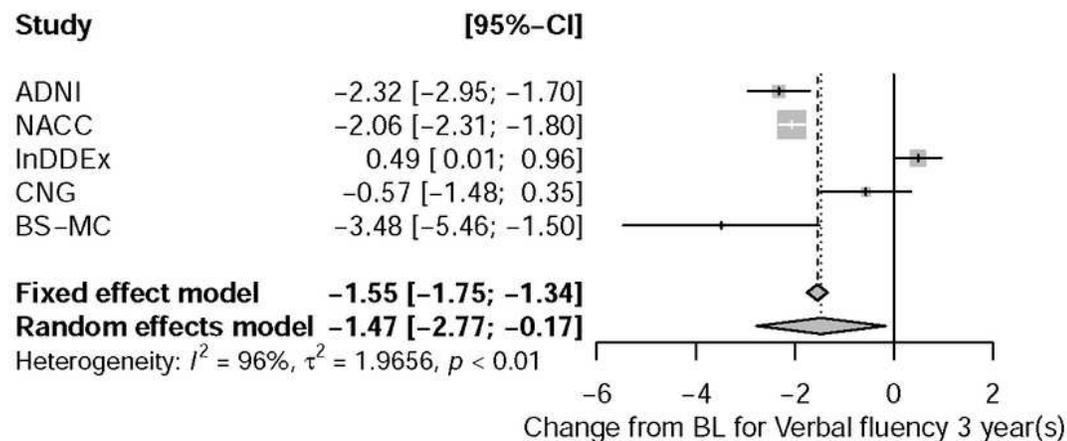
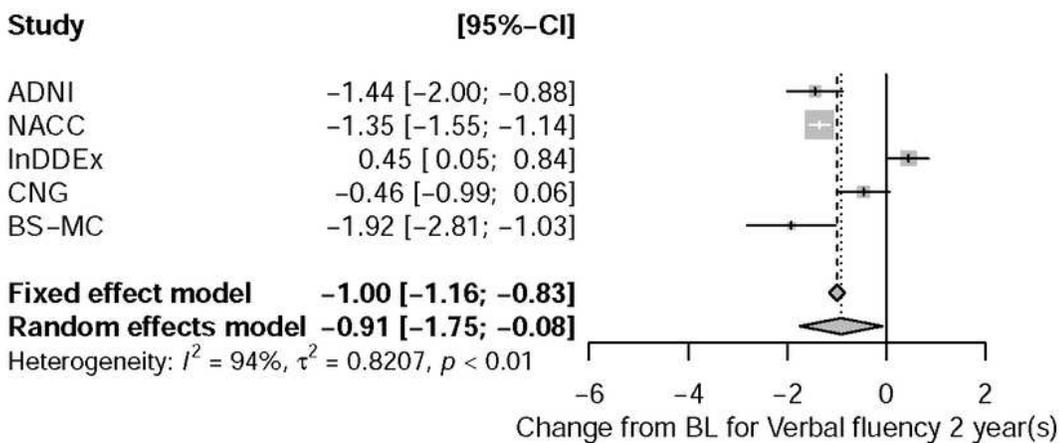
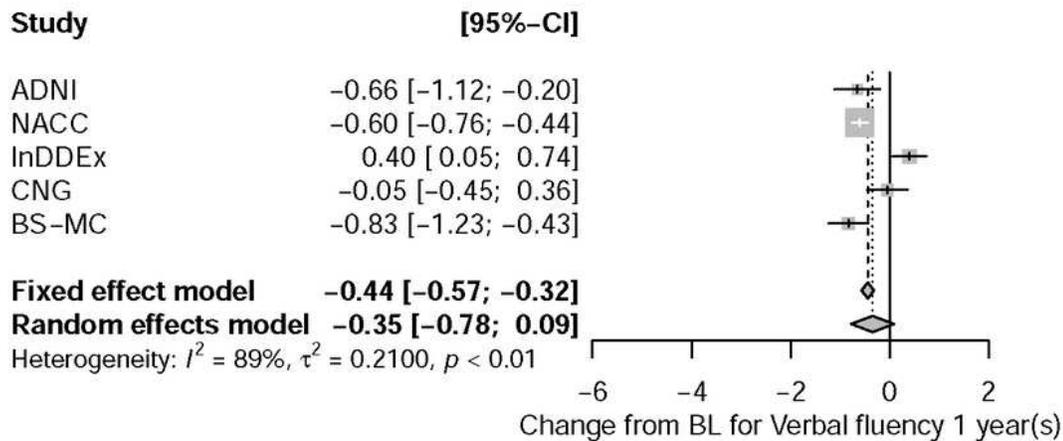


Figure 3

Forest plots for the change of Verbal Fluency (animals) from baseline to 1, 2 and 3 years. Legend: Mean changes and 95% confidence intervals for each study and for the overall effect in the fixed effects and the random effects model are given. τ^2 is the between-study variance, I^2 measures heterogeneity (between study variance over total variance), p -value for the test of heterogeneity. Graphs show study

specific means and confidence intervals for each study as grey squares and lines and for overall effects as diamonds. Size of squares represent precision of individual treatment estimates.

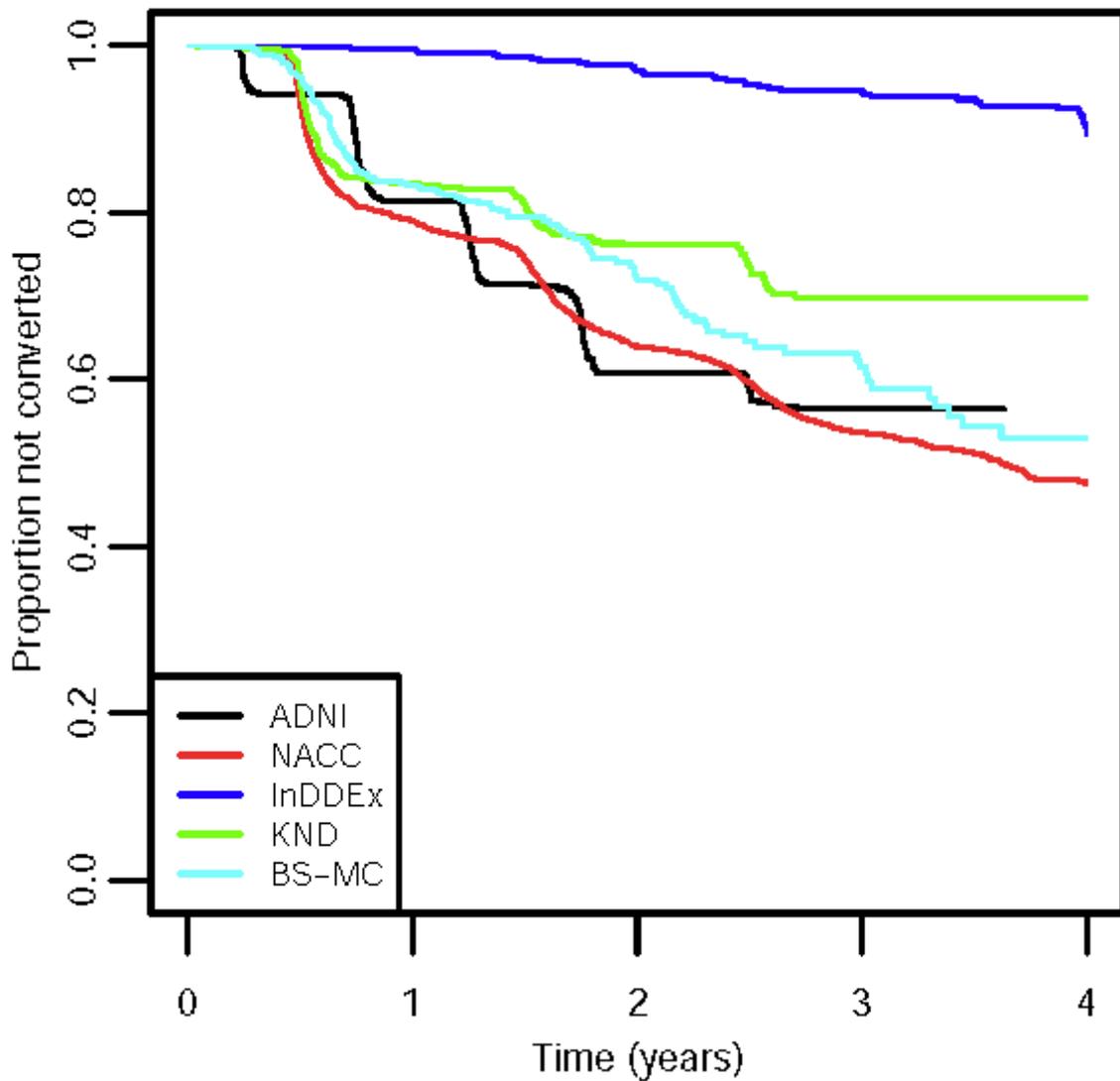


Figure 4

Kaplan-Meier plots of the proportion converted from MCI to dementia versus time for each study.

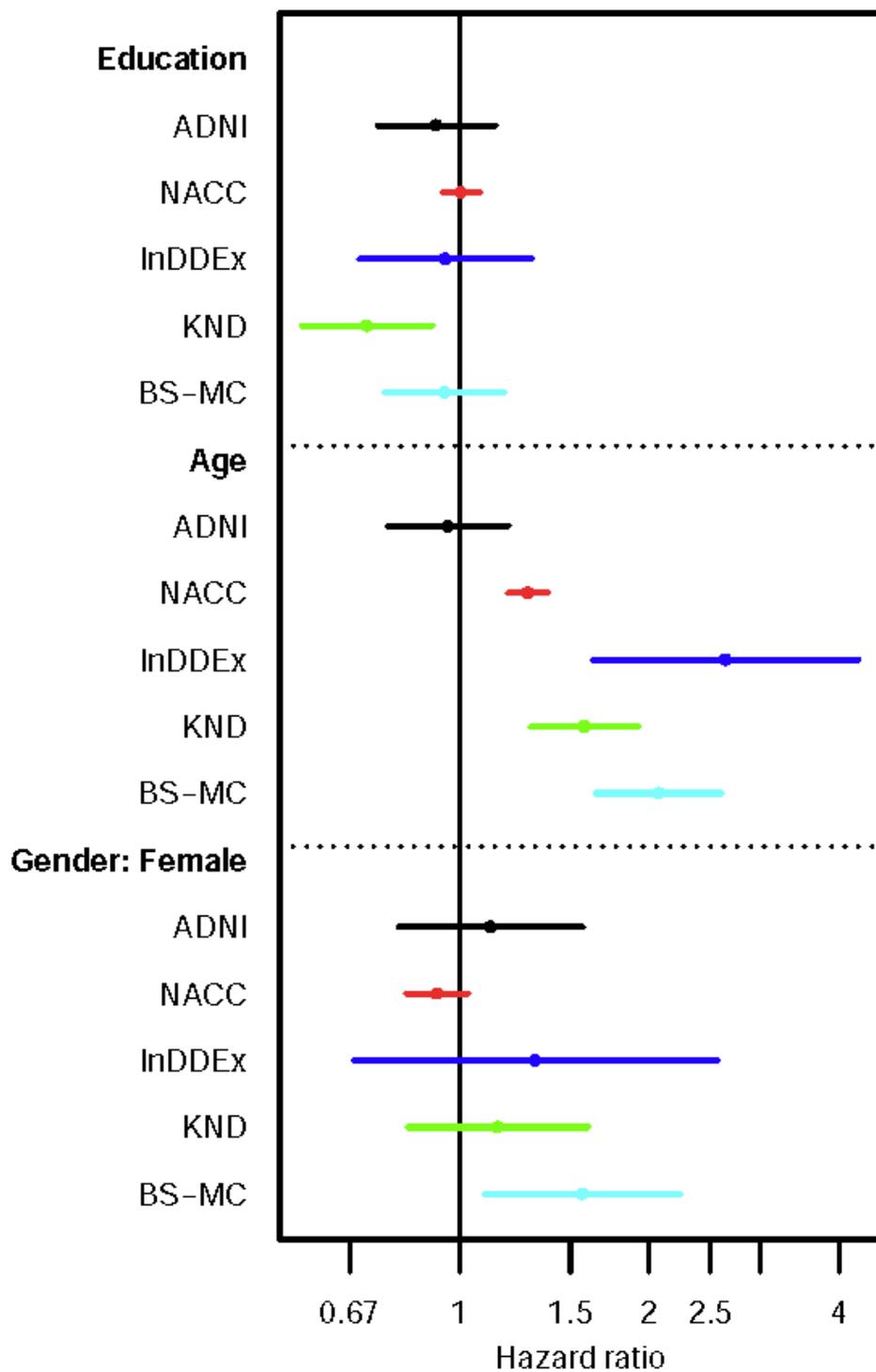


Figure 5

Hazard ratios of five Cox proportional hazard regression models for each study. Legend to Figure 5: For education the hazard ratio for progression to dementia is shown for an increase of 4 years, for age it is shown for an increase of 10 years, for gender it is for females relative to males.