1  **Title**: Development and validation of a harmonized memory score for multicenter Alzheimer's disease

2  and related dementia research

3

4  **Authors**: Mark Sanderson-Cimino[1], Alden L. Gross[2], Leslie S. Gaynor[3,4], Emily W. Paolillo[1], Rowan

5  Saloner[1], Marilyn S. Albert[5], fLiana G. Apostolova[6], Brooke Boersema[7], Adam L. Boxer[1], Bradley F.

6  Boeve[7], Kaitlin B. Casaletto[1,] Savannah R. Hallgarth[1], Valentina E. Diaz[1], Lindsay R. Clark[8], Pauline

7  Maillard[11], Ani Eloyan[9], Sarah Tomaszewski Farias[11], Mitzi M. Gonzales[12], Dustin B. Hammers[6], Renaud

8  La Joie[1], Yann Cobigo[1], Amy Wolf[1], Benjamin M. Hampstead[13], Dawn Mechanic-Hamilton[14], Bruce L.

9  Miller[1], Gil D. Rabinovici[1], John M. Ringman[15], Howie J. Rosen[1], Sephira G. Ryman[16], Jillian L.

10  Prestopnik[16], David P. Salmon[17], Glenn E. Smith[18,19], Charles DeCarli[11], Kumar B. Rajan[20], Lee-Way

11  Jin[21], Jason Hinman[22], David K. Johnson[11], Danielle Harvey[24], Myriam Fornage[25] Joel H. Kramer[1], Adam

12  M. Staffaroni[1], on behalf of the ALLFTD consortium, MarkVCID study, LEADS consortium, and Diverse

13  Vascular Contributions to Cognitive Impairment and Dementia (Diverse VCID) Study Investigators*

14

15  **Affiliations:**
16  [1] Memory and Aging Center, UCSF Weill Institute for Neurosciences, University of California San
17  Francisco, San Francisco, CA, 94158, USA
18  [2] Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD,
19  21205, USA
20  [3] Division of Geriatric Medicine, Department of Medicine, Vanderbilt University Medical Center,
21  Nashville, TN, 37232, USA
22  [4] Vanderbilt Memory and Aging Center, Vanderbilt University Medical Center, Nashville, TN, 37232,
23  USA
24  [5] Department of Neurology, Johns Hopkins University School of Medicine, Baltimore, MD, 21205, USA
25  [6] Department of Neurology, Indiana University School of Medicine, Indianapolis, IN, 46202, USA
26  [7] Department of Neurology, Mayo Clinic, Rochester, MN, 55905, USA
27  [8] Wisconsin Alzheimer's Disease Research Center, University of Wisconsin School of Medicine and
28  Public Health, Madison, WI, 53726, USA
29  [9] Department of Biostatistics, Center for Statistical Sciences, Brown University, Providence, RI, 02912,
30  USA
31  [10] Davis Department of Neurology, University of California, Sacramento, CA, 95816, USA
32  [11] Department of Neurology, University of California at Davis, Sacramento, CA, 95816, USA.
33  [12] Department of Neurology, Cedars Sinai Medical Center, Los Angeles, CA, 90048, USA
34  [13] Department of Psychiatry, University of Michigan, Ann Arbor, MI, 48109, USA

35  [14] Department of Neurology, Perelman School of Medicine at the University of Pennsylvania,
36  Philadelphia, PA, 19104, USA
37  [15] Department of Neurology, Keck School of Medicine at USC, Los Angeles, CA, 90033, USA
38  [16] Center for Memory & Aging, University of New Mexico, Albuquerque, NM, 87110, USA
39  [17] Department of Neurosciences, University of California, San Diego, La Jolla, CA, 92161, USA
40  [18] 1Florida Alzheimer's Disease Research Center, Gainesville, FL, 32610, USA
41  [19] Department of Clinical and Health Psychology, University of Florida, Gainesville, FL, 32603, USA
42  [20] Rush Institute for Healthy Aging, Rush University Medical Center, Chicago, IL, 60612, USA
43  [21] Department of Pathology and Laboratory Medicine University of California, Davis, CA, 95817, USA
44  [22] Department of Neurology, University of California, Los Angeles, Los Angeles, CA, 90095, United
45  States
46  [24] Department of Public Health Sciences University of California, Davis, CA, 95616, USA
47  [25] Brown Foundation Institute of Molecular Medicine, McGovern Medical School at The University of
48  Texas Health Science Center at Houston, Houston, TX, 77030, USA
49
50


51  **Corresponding Authors**:
52          Mark Sanderson-Cimino, PhD
53                  Department of Neurology
54                  675 Nelson Rising Lane, Suite 190
55                  San Francisco, Ca, 94143
56                  415.502.7201
57                  mark.sandersoncimino@ucsf.edu
58
59          Adam Staffaroni, PhD
60                  Department of Neurology
61                  675 Nelson Rising Lane, Suite 190
62                  San Francisco, Ca, 94143
63                  415.502.7201
64                  adam.staffaroni@ucsf.edu
65
66
67


68  **Author Approval:** Authors have approved this manuscript.
69
70
71

72    **Abstract:**

73    INTRODUCTION List-learning tasks are important for characterizing memory in ADRD

74    research, but the Uniform Data Set neuropsychological battery (UDS-NB) lacks a list-learning

75    paradigm; thus, sites administer a range of tests. We developed a harmonized memory composite

76    that incorporates UDS memory tests and multiple list-learning tasks.

77    METHODS: Item-banking confirmatory factor analysis was applied to develop a memory

78    composite in a diagnostically heterogenous sample (n=5943) who completed the UDS-NB and

79    one of five list-learning tasks. Construct validity was evaluated through associations with

80    demographics, disease severity, cognitive tasks, brain volume, and plasma phosphorylated tau (p-

81    tau181 and p-tau217). Test-retest reliability was assessed. Analyses were replicated in a

82    racially/ethnically diverse cohort (n=1058).

83    RESULTS: Fit indices, loadings, distributions, and test-retest reliability were adequate. Expected

84    associations with demographics and clinical measures within development and validation cohorts

85    supported validity.

86    DISCUSSION: This composite enables researchers to incorporate multiple list-learning tasks

87    with other UDS measures to create a single metric.

88

89

## 1. Background

Alzheimer's Disease (AD) and AD-Related Dementia (ADRD) research has been accelerated by the standardization of data collection procedures across multicenter studies. These datasets include larger numbers of participants sampled across geographic regions, leading to more diverse cohorts, greater analytic power, and enhanced generalizability. In the United States, the National Alzheimer's Coordinating Center (NACC) has led a large-scale effort to standardize the collection of common data elements, called the Uniform Data Set (UDS), across more than 35 Alzheimer's Disease Centers (ADC). The UDS includes a detailed collection of data, such as demographic characterization, clinical scales, self- and partner-reported questionnaires, and importantly for AD/ADRD, a battery of neuropsychological assessments [1].

A core element of clinical phenotyping in AD/ADRD is the accurate and robust measurement of memory. The third version of the UDS (UDS v3.0) neuropsychological battery was released in 2015 [1] and includes eleven cognitive tasks, of which only two primarily assess episodic memory: Craft Story (verbal) and Benson Figure (visual). A notable absence from this battery is a list-learning task, one of the most common [2] and effective memory assessments used by researchers and clinicians. Compared to other episodic memory paradigms, list-learning tasks can provide a more detailed characterization of memory profiles, inform differential diagnosis, and improve prediction of dementia progression [3-6].

Although there are practical benefits to allowing ADCs to vary in list-learning task (e.g., comparison to site-specific historic data), a notable downside is limited cross-site comparability. One example of a project affected by site differences in list-learning procedures is Diverse Vascular Contributions to Cognitive Impairment and Dementia (DVCID). DVCID is a prospective, observational study of vascular brain injury that is focused on Black/African Americans and Latina/o/Hispanic Americans. Its primary goals are to investigate basic mechanisms of small vessel cerebrovascular injury and their relationships with social determinants of health, particularly among individuals of diverse backgrounds [7]. When that study began, the decision was made to allow sites to select from one of six

115    list-learning tasks so they could continue to collect data in a manner consistent with their historic

116    procedures.

117          The study presented in this manuscript, borne of the desire to develop a harmonized memory

118    metric for the DVCID consortium, sought to leverage best practices in statistical harmonization

119    techniques [8-10] to create a psychometrically robust memory composite that is on the same scale and

120    interpretable regardless of the list-learning task administered. Such a score has utility not only for

121    DVCID, but many studies seeking to include a memory score for cross-AD/ADRD research. This

122    composite will facilitate comparison of data within a study, even if that study changes list-learning tasks

123    over time. In a recent update to the UDS v4.0, studies are required to administer one of two list-learning

124    tasks (both of which are incorporated in the present study) [11]. As a result, sites and studies will differ in

125    which list-learning task they collect. Furthermore, this requirement will introduce differences in list-

126    learning data relative to retrospectively collected UDS v3.0 data. There is therefore a pressing need for a

127    harmonized memory score.

128          To address this issue, we pooled data from 5943 participants recruited across four consortia and

129    19 ADCs to create and validate the UDS Memory plus List-Learning (UDS-M+) score, a harmonized

130    memory score that accommodates any of five different list-learning tasks. We then test the performance of

131    this composite in a large independent sample, the DVCID cohort (n=1058). We hypothesize a harmonized

132    memory factor composite score 1) shows adequate model fit and reliability; 2) shows evidence for

133    comparability regardless of contributing list-learning task; 3) demonstrates convergent and discriminate

134    validity with other independent cognitive tasks; and 4) shows evidence of construct validity through

135    associations with disease severity, brain volume, and plasma phosphorylated tau (pTau) levels.

136    **2.   Methods**

137    *2.1. Participants:*

138     Participants were recruited from 19 ADCs and four AD/ADRD consortia. The score was first

139     developed and tested in a sample of participants (Development cohort; n=5943) recruited through the

140     ARTFL-LEFFTDS Longitudinal Frontotemporal Lobar Degeneration Study (ALLFTD: NCT04363684),

141     Biomarkers for Vascular Contributions to Cognitive Impairment and Dementia (MarkVCID:

142     NCT06284213) and Longitudinal Early-Onset Alzheimer's Disease Study (LEADS: NCT03507257).

143     Additional data were also included from the 1Florida, Wisconsin, and University of Southern California

144     ADCs. The UDS-M+ was then tested in a separate cohort of participants enrolled through 16 centers in

145     the DVCID consortium (Validation cohort; DVCID: n=1058). Participants in both cohorts (Table 1)

146     presented with a range of clinical diagnoses and severity levels (including cognitively unimpaired) that

147     were made in consensus conferences using published criteria [12-25]. A subset of participants within the

148     Development cohort provided follow-up data (n=462). Only cross-sectional data was used from the

149     Development Cohort; baseline and follow-up data were included from the Validation cohort.

150     *Inclusion/Exclusion*: All participants in both cohorts completed at least one UDS memory test or a

151     list-learning task (see Assessments section). Inclusion criteria required English as the primary language to

152     remove variance associated with language that may impact task performance. An associated study is

153     underway to evaluate the inclusion of tasks completed in other languages. Written informed consent was

154     obtained from all participants or their legal representative. All studies were approved by the Institutional

155     Review Boards of the consortia or ADC in accordance with institutional guidelines and the Helsinki

156     Declaration.

157     ***2.2. Measures:***

158     *List-learning tasks.* Participants were administered one of five list-learning tasks: the California

159     Verbal Learning Test-II, Standard Form (CVLT-II) [26], the California Verbal Learning Test-II, Short

160     Form (CVLT-SF) [26]; Hopkins Verbal Learning Test Revised (HVLT) [27]; Consortium to Establish a

161     Registry for Alzheimer's Disease Word List Memory Test (CERAD) [28]; and the Rey Auditory Verbal

162     Learning Test (AVLT) [29]. All tests begin with an immediate recall paradigm in which examiners read a

163 list of words to the participant, who are asked to repeat as many words as possible. This immediate recall

164 paradigm is repeated several times. The number of words on the list and the number of immediate recall

165 trials varies by task. The total number of words recalled across all immediate recall trials was used for

166 composite generation (immediate recall). After a pre-specified delay period, which differs by task,

167 participants are asked to freely recall as many words as possible (delayed recall). Participants are then

168 asked to identify correct stimuli in a multiple-choice format ("yes" vs "no"), and a recognition memory

169 discriminability index, *d'*, is generated by integrating the number of correct responses with false positives

170 [26].

171    *UDS Memory Measures (Linking Items):* Participants across all list-learning groups completed the

172 UDS Neuropsychological Battery [1], which includes two memory tasks: a story memory task (Craft

173 Story) with verbatim immediate and delayed recall scores; and a visual memory task with a free recall and

174 a recognition trial (correct vs incorrect) of a complex figure (Benson Figure). It also includes the

175 Montreal Cognitive Assessment (MoCA), which provides a brief memory task with a list-learning

176 immediate recall and delayed recall. These three tests provide six memory test items, which are common

177 across all participants.

178 ***2.3. Other assessments***:

179    *Additional Cognitive Tasks.* Participants completed additional UDS measures: Verbal Fluency,

180 Trail Making Test A, Trail Making Test B, Number Span Forward, Number Span Backwards, and the

181 Multilingual Naming Test (MINT) [1, 30]. They also completed a brief cognitive screener (MoCA) [31].

182 Tasks of executive functioning were summarized into a single executive composite score, the UDS3-EF

183 [32], which uses similar IRT methods to those used to generate the UDS-M+. A subset of participants

184 completed a separate memory task, the Tablet-based Cognitive Assessment Tool (TabCAT) Favorites test

185 [33], an associative memory task involving both visual and verbal stimuli. The outcome was total correct

186 responses across immediate and delayed recall conditions.

187        *Clinical Dementia Rating Scale (CDR).* Participants' functional impairment was rated using the

188    CDR, a clinician-administered semi-structured interview with the participant and a study partner [34].

189    Clinicians query the following six areas and provide a rating (values: 0, 0.5, 1, 2, 3): memory, judgment

190    and problem solving, community affairs, home and hobbies, orientation, and personal care. The CDR is

191    commonly used in AD/ADRD research and within clinical trials of AD. A weighted algorithm combines

192    scores within all domains into a global score (CDR-g; scores: 0, 0.5, 1, 2, 3) and sum of boxes (CDR-sb;

193    range: 0-18), with higher scores indicating greater impairment.

194        *Plasma biomarkers:* A subset of Development cohort participants (n=302) provided blood

195    samples at the University of California, San Francisco. Plasma phosphorylated tau-217 (p-tau217) was

196    measured with electrochemiluminescence-based assays on the Meso Scale Discovery platform (MSD,

197    Rockville, MD, USA) using previously published methods [35].

198        A subset of DVCID participants (n=413) provided blood samples that were processed through

199    DVCID, and the Quanterix Single Molecule Array (SiMoA) assay was used to quantify p-tau181. These

200    methods are fully described elsewhere [7].

201        *Imaging:* A sample of participants within the Development cohort (MarkVCID; n=385) and

202    Validation cohort (DVCID; n=856) underwent brain MRI with the same acquisition protocol [7]. MRI

203    acquisition and processing is fully described online (https://markvcid.partners.org/markvcid1-protocols-

204    resources) and in prior publications [36]. Imaging correlates of the memory composite were assessed

205    using a priori regions of interests in subsamples with available processed volumetric MRI data. We

206    created a medial temporal lobe ROI comprising bilateral hippocampal, entorhinal, and parahippocampal

207    regions [37-39]. An occipital lobe ROI was created as a control region [40].

208        To complement the ROI analyses, we leveraged an unbiased whole brain voxel-based

209    morphometry (VBM) approach in a sample of participants at the University of California, San Francisco

210    (n=829) who underwent brain MRI within 90 days of completing the UDS-M+ tasks. These participants

211    were in the Developmental cohort and were distinct from the 385 MarkVCID participants in the imaging

212    analyses described above. Structural MRI data (T1) was acquired with a 3T scanner and processed using

213    Statistical Parametric Mapping in MATLAB (MathWorks, Natick, MA, USA) to conduct a Voxel Based

214    Morphometry (VBM) analysis as previously described [35, 41].

215    *2.4. Confirmatory factor analysis with item banking approach for statistical cocalibration*:

216        *Model Building:* An item-banking confirmatory factor analysis (CFA) approach to harmonization

217    was conducted as described in detail by Gross et al. 2023 and Vonk et al 2022 [10, 42]. Prior to model

218    creation, continuous variables were recoded into ordinal scores with up to 12 response categories and at

219    least twenty observations in each category to facilitate use of a graded response model approach [43]. The

220    equal interval approach we used to recode variables retains the shape of the distribution of the continuous

221    data. In this study, we divided the sample into five groups based on which list-learning task participants

222    completed: CVLT-II, CVLT-SF, CERAD, HVLT, AVLT. All groups also had the six linking items from the

223    aforementioned UDS Neuropsychological Battery. We then fit a CFA model to one of the subsamples to

224    estimate item thresholds and loadings, including both the linking items and test items specific to that

225    subsample. For the current study, an initial model was fit within the list-learning subsample that had the

226    largest sample size (AVLT). This sample was chosen as its large size may reflect a broader range of

227    ability, although prior research has shown that varying the order of models does not meaningfully affect

228    the quality of factor scores produced through an item-banking CFA [10]. This model resulted in item

229    parameters for the three AVLT variables (immediate recall, delayed recall, and recognition) and all six

230    linking items. A subsequent model was then fit to another subset of data (i.e., a new list-learning task), in

231    which the thresholds and loadings for the linking items—UDS memory tasks in this study—were fixed to

232    their corresponding values from the initial model. Thresholds and loadings for variables from the second

233    list-learning task were freely estimated. This procedure was repeated for all subsamples of the data until

234    parameters were estimated for all list-learning variables. Factor scores from this final model, in which all

235    parameters were fixed, using all available participant data were then estimated to derive a single memory

236     factor score that is on the same scale for each participant regardless of which list-learning task was

237     administered. Factor score creation utilized a Maximum likelihood Robust (MLR) estimator and theta

238     parametrization.

239          *Model Fit:* The fit of each model was established separately according to factor loadings (>.4

240     retained [44, 45]), the Tucker Lewis Index (TLI; >0.95 [46]), the Comparative Fit Index (CFI; >0.95

241     [47]), and the Root Mean Square Error of Approximation (RMSEA; <0.05) [48]. Model building was

242     completed in Mplus  [49] via Stata Statistical software  [50] using the RUNMPLUS Stata package  [51] .

243     All other statistical analyses (e.g., regressions, descriptives) were conducted in R (V4.2.3) [52].

244     *2.5. Statistical Analyses*

245          To test for differences between the Development and Validation cohorts on key participant

246     characteristics, linear regressions were fit to compare cohorts on demographic variables, functional

247     impairment, diagnoses, and scores on UDS-M+ component measures (e.g., Craft Story).

248          *Comparing scores across list-learning tasks*: A primary goal of this project was to create a

249     memory composite that was comparable across participants, regardless of which list-learning task was

250     completed. The average UDS-M+ memory composite in the Development cohort was compared across

251     list-learning samples, after adjusting for differences in sample composition, including CDR-g score, age,

252     education, and sex. Cohen's *d* effect sizes were calculated for all pairwise comparisons of UDS-M+

253     memory composites across list-learning groups. Cohen's *d* magnitudes were categorized as small ($\leq 0.20$),

254     medium (0.21 to 0.49), and large ($\geq 0.50$) [53].

255          *Test-Retest reliability:* The correlation between baseline and a second visit (retest interval: mean

256     = 12.67 months; SD=1.4; range: 7-18 months) was used to determine test-retest reliability. This analysis

257     was completed in a subset of participants (n=462) who completed the same list-learning task at baseline

258     and follow-up, as well as met at least one of the following conditions: CDR-g=0, MoCA total score >26,

259     or clinical diagnosis of cognitively unimpaired. These sample restrictions were applied to limit the

260    analysis to those least likely to experience disease-related decline while still maximizing input from each

261    list-learning group and study site.

262         *Marginal reliability*: To examine measurement error of the UDSM+ memory composite across

263    estimated levels of memory ability, marginal reliability was calculated as one minus the square of the

264    standard error of measurement for each participant's UDS-M+ score [54]. This metric was then plotted

265    against the latent factor to create a test information plot that displays reliability of the latent factor across

266    the latent trait, as measured by the UDS-M+. In general, reliability above .80 is recommended [10]. Test

267    information plots were created for the full sample, then delineated by list-learning group.

268    *2.6. Validity analyses*:

269         To evaluate evidence for convergent construct validity, separate linear regressions were fit to test

270    the association between the UDS-M+ and each of the following: age, sex, education, and TabCAT

271    Favorites. Each regression adjusted for the CDR-sb, unless CDR-sb was the primary variable of interest.

272    Known-groups validity was tested by fitting linear regressions to examine the effect of CDR-g on the

273    memory composite score. Post-hoc linear regressions investigated the relationship between the UDS-M+

274    and the CDR memory box score. Discriminant construct validity was evaluated by testing associations

275    with non-memory cognitive tasks [32]. Analyses were first completed in the Development cohort then

276    replicated as possible within the Validation cohort.

277         Given that memory deficits are an early and common consequence of AD pathology, the

278    association between UDS-M+ score and log (10)-transformed pTau levels were tested via regressions that

279    controlled for age, sex, and education.

280         The relationship between UDS-M+ and imaging (MTL, occipital) was examined using linear

281    regressions controlling for age, total intracranial volume, education, and CDR-sb. The VBM completed in

282    a subset of the Development cohort included age, total intracranial volume, education, and CDR-sb as

283    covariates. A family wise-error threshold was applied to VBM results (Monte-Carlo; 1000 permutations)

284    to generate an error distribution with a cut off at the Type 1 error threshold (95%). This distribution yields

285    a critical T threshold at a family wise error threshold of p<.05 [55].

**3. Results:**

287          *Sample Characteristics:* Participant characteristics are provided in Table 1. Compared to the

288    Development cohort, Validation participants were significantly older and less impaired on the MoCA and

289    CDR. Sample characteristics delineated by list-learning task are presented in Supplemental Table 1.

290          *Model Building and Fit*: In the initial model, all linking item standardized loadings were

291    acceptable (range: 0.60 to 0.93). Linking item loadings were carried forward to all other list-learning

292    models. Fit for all models was adequate or excellent for most fit statistics (CFIs > 0.99: TLIs > 0.99;

293    RMSEA range: 0.03 to 0.09). Fit statistics are summarized in Supplemental Table 2. Unstandardized

294    loadings are presented in Supplemental Table 3 and thresholds for all UDS-M+ variables are presented in

295    Supplemental Table 4.

296          *Psychometrics and cross-sample comparisons*: The UDS-M+ factor score was normally

297    distributed within the full sample and within each list-learning group (Figure 1). After adjusting for

298    differences in sample characteristics (i.e., age, sex, education, and CDR-sb) only the CVLT-II group

299    demonstrated an average Cohen's *d* value in the medium range (average Cohen's *d* = -.40). The average

300    Cohen's *d* estimates were minimal for other list-learning groups (range: -0.05 to -0.24). In participants

301    with longitudinal data (n=462) who were categorized as unimpaired (CDR-g = 0, MoCA > 26, and/or

302    cognitively unimpaired diagnosis), test-retest reliability was adequate (r=0.67; 95% CI: 0.62, 0.72;

303    p<0.001). Marginal reliability was above 0.8 for most participants (93%), only dropping below 0.8 at the

304    highest and lowest ability levels (Figure 2). Supplemental Figure 1 presents marginal reliability delineated

305    by CDR-g groups.

306          *Associations with Clinical Measures:* Associations between the UDS-M+ and demographics and

307    clinical measures are displayed as standardized regression betas in Figure 3. Within the Development

308    cohort higher UDS-M+ scores were associated as expected with younger age ($\beta$=-0.11; 95% CI: -0.09, -

309    0.13; p<0.001), greater years of education ($\beta$=0.12; 95% CI: 0.10, 0.14; p<0.001), and female sex

310    ($\beta$=0.17; 95% CI: 0.13, 0.21; p<0.001). Lower scores were associated with greater functional impairment,

311    as measured by the CDR-sb ($\beta$=-0.65; 95% CI: -0.67, -0.62; p<0.001) and worse performance on a

312    cognitive screener (MoCA: $\beta$=-0.62; 95% CI: -0.60, -0.64; p<0.001). The UDS-M+ was strongly

313    correlated with an independent memory test (TabCAT Favorites: $\beta$=0.59; 95% CI: 0.53, 0.66; p<0.001),

314    whereas the magnitude of associations with measures of executive functioning ($\beta$ range: 0.17, 0.45) and

315    language ($\beta$ range: 0.21, 0.40) were relatively lower. There was a stepwise decline (Figure 4) in UDS-M+

316    performance across CDR-g scores (CDR 0 > 0.5 > 1 > 2+; p's <0.001) and memory box scores (0>1>2+;

317    p<0.001). As shown in Figures 3 and 4, the UDS-M+ associations showed the same general pattern of

318    results in the Development and Validation cohorts. Two associations were notably different between the

319    cohorts. The association between the UDS-M+ and CDR-sb was lower in the Validation cohort than the

320    Development cohort (r = 0.64; 95% CI: 0.62, 0.67; vs r = 0.28 95% CI: 0.24, 0.33). Similarly, the UDS-

321    M+ was less associated with MoCA global scores in the Validation cohort than in the Development

322    Cohort (r = 0.62; 95% CI: 0.60, 0.65 vs r = 0.43; 95% CI: 0.39, 0.46). As noted in Table 1, the cohorts

323    differ in the average and distribution of CDR-sb and MoCA scores, with DVCID representing a less

324    impaired or symptomatic group with restricted variance.

325        *Associations with pTau and brain volume*:  Worse performance on the UDS-M+ was strongly

326    associated with higher plasma p-tau217 levels in the Development cohort (Figure 5A), after controlling

327    for age, education, sex, and CDR-sb (*$\beta$*=-0.37; 95% CI: -0.46, -0.27; p<0.001). Similarly, in DVCID

328    (Figure 5B), worse performance on the UDS-M+ was associated with higher plasma levels of p-tau181 ($\beta$

329    = -0.15; 95% CI: -0.24, -0.05; p=0.002), with the strongest association in the CDR-g = 0.5 group ($\beta$ = -

330    0.33; 95% CI: -0.49, -0.17; p<0.001).

331        Lower UDS-M+ scores were associated with smaller medial temporal lobe volumes in the

332    Development cohort ($\beta$ = 0.23 95% CI: 0.13, 0.33; p<.001) and DVCID ($\beta$ = 0.11 95% CI: 0.04, 0.17;

333    p=0.001). Regarding divergent validity, the UDS-M+ was not significantly associated with the occipital

334    region ($\beta$ = 0.01 95% CI: -0.05, 0.07; p=.68). In a separate subsample, using a whole-brain VBM

335    approach, worse UDS-M+ scores were associated with smaller brain volumes, with the largest cluster

336    including the medial temporal lobe (Supplemental table 5 and Supplemental figure 2), particularly on the

337    left side of the brain.

**4. Discussion:**

339        We present the development and validation of the UDS-M+, a memory composite score

340    developed using advanced psychometric methods  [10, 56] to enable multi-cohort AD/ADRD research.

341    The UDS-M+ was built in a diagnostically heterogeneous sample to support generalizability, and the

342    resulting score showed good model fit with moderate-to-strong factor loadings, and good psychometric

343    properties including a normal distribution and overlapping distributions regardless of the list-learning task

344    administered. Test-retest reliability was adequate and similar to that reported for list-learning tasks over

345    similar time frames [57, 58]. Construct validity was supported by strong associations in the hypothesized

346    directions with demographic factors (e.g., age), an independent memory task, medial temporal lobe

347    volume, and plasma levels of p-tau. The robustness of these findings is bolstered by replication within an

348    independent Validation cohort.

349        The UDS-M+ was developed using best practices in statistical harmonization, allowing

350    researchers to combine the UDS v3.0 neuropsychological tests and five distinct list-learning tasks into a

351    single composite score. Our results, particularly the minimal residual differences between list-learning

352    task groups (Figure 1) and test information plots (Figure 2), suggest the harmonization process was

353    successful. Importantly, the success of this harmonization is partly demonstrated by the observation that

354    some *differences* in memory composite scores were retained across list-learning tasks and studies. First,

355    the list-learning tasks vary in their length, delay interval, and stimuli. These differences impact the

356    difficulty of the task (i.e., CVLT-II difficulty > CVLT-SF) which in turn affects the relationship between

357    the list-learning task and the UDS-M+. As a result, list-leaning tasks vary somewhat in factor loadings

358   (Supplemental table 3), average UDS-M+ scores (Figure 1), and reliability across levels of memory

359   ability (Figure 2). These variations likely reflect real differences across list-learning tasks and support the

360   use of an IRT-based approach over more simplistic composite approaches (e.g., Z score average) that treat

361   all tasks as equal assessments of the construct. Second, the list-learning tasks were not assigned

362   randomly; for example, nearly all members of a healthy aging study (UCSF: Brain Aging Network for

363   Cognitive Health; BRANCH) completed the CVLT-II Standard Form while all participants of an EOAD

364   study (LEADS) completed the AVLT. A comparison of UDS-M+ scores across these list-learning groups

365   reveals, as expected, a difference in means (Figure 1; CVLT-SF > AVLT) and item reliability across levels

366   of memory (Figure 2).

367          An important feature of the UDS-M+ is that the resulting memory score is on the same scale

368   regardless of which list-learning task was administered, This feature reflects a significant improvement

369   over the use of raw scores or standardized scores, which are not directly comparable across samples [59-

370   61].  For example, say participant A has a CVLT-II delayed recall of nine words and participant B has a

371   CVLT-SF delayed recall of nine words. These scores do not reflect equal memory abilities despite

372   recalling the same number of words, as there are notable differences in task demands (e.g., word-list

373   length, recall interval, stimuli). An alternative harmonization approach is to convert scores on each list-

374   learning task to Z scores, based on the distribution of the list-learning task within each sample. A

375   complication with Z score comparisons is that the list-learning tasks likely differ in distribution and, like

376   raw scores, do not account for differences in task demands. For example, if participant A (CVLT-II) and

377   participant B (CVLT-SF) both have a Z score of 0 it does not necessarily mean that they have equal

378   memory ability, just that they have average scores within their respective subsamples. The UDS-M+, in

379   contrast, was created with IRT methods that allow for differences in task demands and truly place each

380   participant's score on the same scale. As a result, if participant A (CVLT-II) had a UDS-M+ score of 0 and

381   participant B (CVLT-SF) had a UDS-M+ score of 0, we can conclude that participant A's memory is

382   measured to be equal to that of participant B (with a margin of error).

383    The ability to study memory regardless of list-learning task will become particularly important as

384    the UDS v4.0 requires investigators to administer either the AVLT or CERAD [11]. Thus, many UDS sites

385    who have historically administered a different list-learning task will face a dilemma: they may switch to a

386    new task or administer two list-learning tasks. Switching to a new task disrupts alignment with all

387    retrospectively collected data. Administering two tasks is burdensome for participants and staff and

388    potentially introduces interference effects. This memory composite may offer a solution, as the UDS-M+

389    can be used to harmonize the AVLT and CERAD with other tasks.

390    The development of this memory composite adds to several ongoing efforts to advance cognitive

391    research via best practices in statistical harmonization. This item-banking approach has been used to

392    harmonize cognitive instruments administered in six countries through the Health and Retirement Study's

393    Harmonized Cognitive Assessment Protocol (HRS-HCAP) and has been applied to other longitudinal

394    cohorts [8, 10, 42, 62-64]. AD/ADRD item-banking studies have shown exceptional methodological

395    validation [8, 10, 42, 62, 64], but only a subset have demonstrated convergent validity or included

396    biomarker validation [61, 65]. Our study included both types of validation and provided replication in an

397    independent cohort with racial/ethnic diversity, further highlighting the value of these modern

398    psychometric techniques to AD/ADRD research. Our study also provides a template for other large

399    studies that seek to either combine somewhat disparate cognitive data—for example, various visuospatial

400    tasks—or who wish to alter their data collection. This later point is especially relevant as there is

401    presently a rapid expansion of digital cognitive assessment tools [66-68]; many long-running studies,

402    including the UDS, are modernizing their cognitive batteries and there is a long-overdue push for

403    inclusion of culturally-appropriate cognitive measures [69, 70].

404    There are limitations to the current version of this score. Marginal reliability was sufficient across

405    all list-learning tasks, though estimates were lower at the high and low extremes of latent memory ability.

406    At least some of the differences in marginal reliability are likely due to sample differences (e.g., disease

407    severity ranges) and task differences (e.g., CVLT-II difficulty > CVLT-SF). As such, it is important for

408     studies that use the UDS-M+ to select list-learning tasks that are appropriate to their sample and,

409     potentially, that vary across participants. For example, the UDS-M+ may better harmonize within a study

410     if impaired participants receive the CVLT-SF and unimpaired participants complete the CVLT-II.

411     Additionally, we primarily present cross-sectional data and have not systematically investigated

412     longitudinal change on the UDS-M+. Ongoing longitudinal analyses will be presented in a separate

413     publication. Although results from the Development cohort support the use of the score across a range of

414     disease severities and diagnoses, and replication in DVCID support its use with participants identifying as

415     Black/African American and Latino/a/Hispanic Americans, additional studies are required to extend the

416     UDS-M+ for use in other testing languages and cultures. Work is underway to develop a version of the

417     score for Spanish speakers. The UDS-M+ was created to aid in comparing memory scores across

418     participants but should not be interpreted as a normatively corrected estimate of memory ability. A score

419     of 0, for example, is relative to this mixed diagnostic group, and should therefore not be interpreted as

420     representing average memory ability at a population level. We also do not provide demographic

421     adjustments, and recommend this step be conducted by investigators in their statistical models based on

422     the relevant confounds in their study. Although the score can accommodate five of the most common list-

423     learning tasks, there are other measures that were not included. The code used to generate the UDS-M+,

424     however, is designed to be iterative, allowing for the addition of other list-learning tasks and memory

425     measures in the future. Interested users with UDS v2.0 data could consider using the crosswalk study to

426     convert Logical Memory to Craft Story scores, although validation is strongly suggested. The code to

427     create the UDS-M+ will be made publicly available at the time of publication.

428     **Conclusion**

429          We have successfully applied item-banking methods to harmonize data across and within ADCs.

430     The UDS-M+ provides each participant a memory composite score that integrates information across

431     multiple memory measures, regardless of which list-learning task was completed. Without the UDS-M+,

432     ADCs and multisite consortia would be forced to either ignore list-learning data or to combine memory

433    scores in a suboptimal manor [61, 71]. The UDS-M+ may therefore substantially expand and improve the

434    study of memory within ADCs. The successful creation of the UDS-M+ serves as a model for future

435    studies wishing to harmonize internally or who plan to alter their neuropsychological battery during

436    conversion to the UDS v4.0.

437

438                                        References

439    1.    Weintraub, S., et al., *Version 3 of the Alzheimer Disease Centers' neuropsychological test battery*
440          *in the Uniform Data Set (UDS).* Alzheimer disease and associated disorders, 2018. **32**(1): p. 10.
441    2.    Rabin, L.A., E. Paolillo, and W.B. Barr, *Stability in test-usage practices of clinical*
442          *neuropsychologists in the United States and Canada over a 10-year period: A follow-up survey of*
443          *INS and NAN members.* Archives of Clinical Neuropsychology, 2016. **31**(3): p. 206-230.
444    3.    Mansbach, W.E., R.A. Mace, and K.M. Clark, *Story recall and word lists: Differential and*
445          *combined utilities in predicting cognitive diagnosis.* Journal of Clinical and Experimental
446          Neuropsychology, 2014. **36**(6): p. 569-576.
447    4.    De Simone, M.S., et al., *Different deficit patterns on word lists and short stories predict*
448          *conversion to Alzheimer's disease in patients with amnestic mild cognitive impairment.* Journal of
449          Neurology, 2017. **264**: p. 2258-2267.
450    5.    Tremont, G., et al., *Comparison of verbal memory impairment rates in mild cognitive impairment.*
451          Journal of Clinical and Experimental Neuropsychology, 2010. **32**(6): p. 630-636.
452    6.    Tremont, G., et al., *Differential impact of executive dysfunction on verbal list learning and story*
453          *recall.* The Clinical Neuropsychologist, 2000. **14**(3): p. 295-302.
454    7.    DeCarli, C., et al., *WMH Contributions to Cognitive Impairment: Rationale and Design of the*
455          *Diverse VCID Study.* Stroke, 2024.
456    8.    Kobayashi, L.C., et al., *Cross national comparisons of later life cognitive function using data*
457          *from the Harmonized Cognitive Assessment Protocol (HCAP): Considerations and recommended*
458          *best practices.* Alzheimer's & Dementia, 2024. **20**(3): p. 2273-2281.
459    9.    Briceño, E.M., et al., *A cultural neuropsychological approach to harmonization of cognitive data*
460          *across culturally and linguistically diverse older adult populations.* Neuropsychology, 2023.
461          **37**(3): p. 247.
462    10.   Gross, A.L., et al., *Harmonisation of later-life cognitive function across national contexts: results*
463          *from the Harmonized Cognitive Assessment Protocols.* The Lancet Healthy Longevity, 2023.
464          **4**(10): p. e573-e583.
465    11.   Center, A.C.T.F.a.N.A.s.C., *Instructions For the Neuropsychological Battery (Form C2).* 2025.
466    12.   Petersen, R.C., et al., *Alzheimer's disease neuroimaging initiative (ADNI): clinical*
467          *characterization.* Neurology, 2010. **74**(3): p. 201-209.
468    13.   Petersen, R.C., et al., *Mild cognitive impairment: a concept in evolution.* Journal of internal
469          medicine, 2014. **275**(3): p. 214-228.
470    14.   Dubois, B., et al., *Research criteria for the diagnosis of Alzheimer's disease: revising the*
471          *NINCDS–ADRDA criteria.* The Lancet Neurology, 2007. **6**(8): p. 734-746.
472    15.   Rascovsky, K., et al., *Sensitivity of revised diagnostic criteria for the behavioural variant of*
473          *frontotemporal dementia.* Brain, 2011. **134**(9): p. 2456-2477.
474    16.   Neary, D., et al., *Frontotemporal lobar degeneration: a consensus on clinical diagnostic criteria.*
475          Neurology, 1998. **51**(6): p. 1546-1554.
476    17.   Postuma, R.B., et al., *MDS clinical diagnostic criteria for Parkinson's disease.* Movement
477          disorders, 2015. **30**(12): p. 1591-1601.
478    18.   Sachdev, P., et al., *Diagnostic criteria for vascular cognitive disorders: a VASCOG statement.*
479          Alzheimer Disease & Associated Disorders, 2014. **28**(3): p. 206-218.
480    19.   Litvan, I., et al., *Accuracy of the clinical diagnoses of Lewy body disease, Parkinson disease, and*
481          *dementia with Lewy bodies: a clinicopathologic study.* Archives of neurology, 1998. **55**(7): p.
482          969-978.
483    20.   Gorno-Tempini, M.L., et al., *Classification of primary progressive aphasia and its variants.*
484          Neurology, 2011. **76**(11): p. 1006-1014.
485    21.   McKeith, I.G., et al., *Diagnosis and management of dementia with Lewy bodies: Fourth*
486          *consensus report of the DLB Consortium.* Neurology, 2017. **89**(1): p. 88-100.

22. Brooks, B.R., et al., *El Escorial revisited: revised criteria for the diagnosis of amyotrophic lateral sclerosis.* Amyotrophic lateral sclerosis and other motor neuron disorders, 2000. **1**(5): p. 293-299.

23. Petersen, R.C., et al., *Mild cognitive impairment due to Alzheimer disease in the community.* Annals of neurology, 2013. **74**(2): p. 199-208.

24. Jack Jr, C.R., et al., *NIA‐AA research framework: toward a biological definition of Alzheimer's disease.* Alzheimer's & Dementia, 2018. **14**(4): p. 535-562.

25. Barker, M.S., et al., *Proposed research criteria for prodromal behavioural variant frontotemporal dementia.* Brain, 2022. **145**(3): p. 1079-1097.

26. Delis, D., et al., *California verbal learning test, adult version (CVLT-II).* Cleveland, Ohio: The Psychological Corporation, 2000.

27. Benedict, R.H., et al., *Hopkins Verbal Learning Test–Revised: Normative data and analysis of inter-form and test-retest reliability.* The Clinical Neuropsychologist, 1998. **12**(1): p. 43-55.

28. Morris, J.C., et al., *The Consortium to Establish a Registry for Alzheimer's Disease (CERAD). Part I. Clinical and neuropsychological assessment of Alzheimer's disease.* Neurology, 1989. **39**(9): p. 1159-1165.

29. Schmidt, M., *Rey auditory verbal learning test: A handbook.* Vol. 17. 1996: Western Psychological Services Los Angeles, CA.

30. Besser, L., et al., *Version 3 of the national Alzheimer's coordinating center's uniform data set.* Alzheimer Disease & Associated Disorders, 2018. **32**(4): p. 351-358.

31. Nasreddine, Z.S., et al., *The Montreal Cognitive Assessment, MoCA: a brief screening tool for mild cognitive impairment.* Journal of the American Geriatrics Society, 2005. **53**(4): p. 695-699.

32. Staffaroni, A.M., et al., *Development and validation of the Uniform Data Set (v3. 0) executive function composite score (UDS3‐EF).* Alzheimer's & dementia, 2021. **17**(4): p. 574-583.

33. Thompson, L.I., et al., *Remote and in‐clinic digital cognitive screening tools outperform the MoCA to distinguish cerebral amyloid status among cognitively healthy older adults.* Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring, 2023. **15**(4): p. e12500.

34. Morris, J.C., *Clinical dementia rating: a reliable and valid diagnostic and staging measure for dementia of the Alzheimer type.* International psychogeriatrics, 1997. **9**(S1): p. 173-176.

35. Saloner, R., et al., *Plasma phosphorylated tau‐217 exhibits sex‐specific prognostication of cognitive decline and brain atrophy in cognitively unimpaired adults.* Alzheimer's & Dementia, 2024. **20**(1): p. 376-387.

36. Lu, H., et al., *MarkVCID cerebral small vessel consortium: II. Neuroimaging protocols.* Alzheimer's & Dementia, 2021. **17**(4): p. 716-725.

37. Aljabar, P., et al., *Multi-atlas based segmentation of brain images: atlas selection and its effect on accuracy.* Neuroimage, 2009. **46**(3): p. 726-38.

38. Desikan, R.S., et al., *An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest.* Neuroimage, 2006. **31**(3): p. 968-980.

39. Fischl, B., et al., *Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain.* Neuron, 2002. **33**(3): p. 341-55.

40. Carmichael, O., et al., *MRI predictors of cognitive change in a diverse and carefully characterized elderly population.* Neurobiol Aging, 2012. **33**(1): p. 83-95.

41. Mandelli, M.L., et al., *Frontal white matter tracts sustaining speech production in primary progressive aphasia.* Journal of Neuroscience, 2014. **34**(29): p. 9754-9767.

42. Vonk, J.M., et al., *Cross-national harmonization of cognitive measures across HRS HCAP (USA) and LASI-DAD (India).* Plos one, 2022. **17**(2): p. e0264166.

43. Samejima, F., *The general graded response model*, in *Handbook of polytomous item response theory models*. 2011, Routledge. p. 87-118.

44. Williams, B., A. Onsman, and T. Brown, *Exploratory factor analysis: A five-step guide for novices.* Australasian journal of paramedicine, 2010. **8**: p. 1-13.

45. Matsunaga, M., *How to Factor-Analyze Your Data Right: Do's, Don'ts, and How-To's.* International journal of psychological research, 2010. **3**(1): p. 97-110.

538  46.  Schumacker, R.E. and R.G. Lomax, *A beginner's guide to structural equation modeling*. 2004:
539        psychology press.
540  47.  Savalei, V. and P.M. Bentler, *Structural equation modeling*. The handbook of marketing research:
541        Uses, misuses, and future advances, 2006. **330**: p. 36.
542  48.  Byrne, B.M., *Structural equation modeling: Perspectives on the present and the future.*
543        International Journal of Testing, 2001. **1**(3-4): p. 327-334.
544  49.  Kelloway, E.K., *Using Mplus for structural equation modeling: A researcher's guide*. 2014: Sage
545        Publications.
546  50.  StataCorp, L., *Stata statistical software: release 15 college station.* TX StataCorp LP, 2023. **5**: p.
547        231-9.
548  51.  Jones, R., *RUNMPLUS: Stata module to run Mplus from Stata.* 2013.
549  52.  R Core Team, R. and R.C. Team, *R Foundation for Statistical Computing; Vienna, Austria: 2020*.
550        2023.
551  53.  Gignac, G.E. and E.T. Szodorai, *Effect size guidelines for individual differences researchers.*
552        Personality and individual differences, 2016. **102**: p. 74-78.
553  54.  Green, B.F., et al., *Technical guidelines for assessing computerized adaptive tests.* Journal of
554        Educational measurement, 1984. **21**(4): p. 347-360.
555  55.  Kimberg, D.Y., H.B. Coslett, and M.F. Schwartz, *Power in voxel-based lesion-symptom mapping.*
556        Journal of cognitive neuroscience, 2007. **19**(7): p. 1067-1080.
557  56.  Crane, P.K., et al., *Measurement precision across cognitive domains in the Alzheimer's Disease*
558        *Neuroimaging Initiative (ADNI) data set.* Neuropsychology, 2023. **37**(4): p. 373.
559  57.  Alioto, A.G., et al., *Long-term test-retest reliability of the California Verbal Learning Test–second*
560        *edition.* The Clinical Neuropsychologist, 2017. **31**(8): p. 1449-1458.
561  58.  Stein, J., et al., *The assessment of changes in cognitive functioning: age-, education-, and gender-*
562        *specific reliable change indices for older adults tested on the CERAD-NP battery: results of the*
563        *German Study on Ageing, Cognition, and Dementia in Primary Care Patients (AgeCoDe).* The
564        American Journal of Geriatric Psychiatry, 2012. **20**(1): p. 84-97.
565  59.  Lacritz, L.H., et al., *Comparison of the hopkins verbal learning test-revised to the California*
566        *verbal learning test in Alzheimer's disease.* Applied Neuropsychology, 2001. **8**(3): p. 180-184.
567  60.  Johnson, C., *Similarities and Differences among Commonly Used Verbal List Learning Tasks.*
568        2012.
569  61.  Hampton, O.L., et al., *Harmonizing the preclinical Alzheimer cognitive composite for multicohort*
570        *studies.* Neuropsychology, 2023. **37**(4): p. 436.
571  62.  Mukherjee, S., et al., *Cognitive domain harmonization and cocalibration in studies of older*
572        *adults.* Neuropsychology, 2022.
573  63.  Yi, D., et al., *The Korean brain aging study for the early diagnosis and prediction of Alzheimer's*
574        *disease (KBASE): Cognitive data harmonization.* Alzheimer's & Dementia, 2023. **19**: p. e064533.
575  64.  Scollard, P., et al., *Ceiling effects and differential measurement precision across calibrated*
576        *cognitive scores in the Framingham Study.* Neuropsychology, 2023. **37**(4): p. 383.
577  65.  Choi, S.E., et al., *Development and validation of language and visuospatial composite scores in*
578        *ADNI.* Alzheimer's & Dementia: Translational Research & Clinical Interventions, 2020. **6**(1): p.
579        e12072.
580  66.  Tsoy, E., S. Zygouris, and K.L. Possin, *Current state of self-administered brief computerized*
581        *cognitive assessments for detection of cognitive disorders in older adults: a systematic review.*
582        The journal of prevention of Alzheimer's disease, 2021. **8**: p. 267-276.
583  67.  Öhman, F., et al., *Current advances in digital cognitive assessment for preclinical Alzheimer's*
584        *disease.* Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring, 2021. **13**(1): p.
585        e12217.
586  68.  Staffaroni, A.M., et al., *Digital Cognitive Assessments for Dementia: Digital assessments may*
587        *enhance the efficiency of evaluations in neurology and other clinics.* Practical Neurology (Fort
588        Washington, Pa.), 2020. **2020**: p. 24.

69. Franzen, S., et al., *Cross-cultural neuropsychological assessment in Europe: Position statement of the European consortium on cross-cultural neuropsychology (ECCroN).* The Clinical Neuropsychologist, 2022. **36**(3): p. 546-557.

70. Merkley, T.L., et al., *Challenges and opportunities for harmonization of cross-cultural neuropsychological data.* Neuropsychology, 2023. **37**(3): p. 237.

71. McNeish, D. and M.G. Wolf, *Thinking twice about sum scores.* Behavior research methods, 2020. **52**: p. 2287-2305.

**Funding**

**Conflicts of interest**

*B.F.B.* has served as an investigator for clinical trials sponsored by Alector, Biogen, Transposon and Cognition Therapeutics. He serves on the Scientific Advisory Board of the Tau Consortium, which is funded by the Rainwater Charitable Foundation. He receives research support from NIH. *A.L.B.* receives research support from the NIH, the Tau Research Consortium, the Association for Frontotemporal Degeneration, Bluefield Project to Cure Frontotemporal Dementia, Corticobasal Degeneration Solutions, the Alzheimer's Drug Discovery Foundation and the Alzheimer's Association. He has served as a consultant for Aeovian, AGTC, Alector, Arkuda, Arvinas, Boehringer Ingelheim, Denali, GSK, Life Edit, Humana, Oligomerix, Oscotec, Roche, TrueBinding and Wave and received research support from Biogen, Eisai and Regeneron. *M.M.G.* reports personal stock in Abbvie. R.C.P. reports personal fees from Roche, no personal fees from Eisai, and personal fees from Genentech, personal fees from Eli Lilly and personal fees from Nestle, outside the submitted work. *B.L.M* reported serving on the scientific advisory board of the Bluefield Project to Cure Frontotemporal Dementia; the John Douglas French Alzheimer's Foundation; Fundación Centro de Investigación Enfermedades Neurológicas, Madrid, Spain; Genworth; the Kissick Family Foundation; the Larry L. Hillblom Foundation; and the Tau Consortium of the Rainwater Charitable Foundation; serving as a scientific advisor for the Arizona Alzheimer's Consortium; Massachusetts General Hospital Alzheimer's Disease Research Center; and the Stanford University Alzheimer's Disease Research Center; receiving royalties from Cambridge University Press, Elsevier, Guilford Publications, Johns Hopkins Press, Oxford University Press, and the Taylor & Francis Group; serving as editor for Neurocase and section editor for Frontiers in Neurology; and receiving grants for the University of California San Francisco Frontotemporal Dementia Core, from the Bluefield Project to Cure Frontotemporal Dementia, and from the National Institute on Aging for the US–South American Initiative for Genetic-Neural-Behavioral Interactions in Human Neurodegenerative Diseases. *G.D.R.* reported grants from National Institutes of Health during the conduct of the study; consulting fees from C2N, Eli Lilly, Alector, Merck, Roche, and Novo Nordisk; data safety monitoring board fees from Johnson & Johnson; and grants from Avid Radiopharmaceuticals, GE Healthcare, Life Molecular Imaging, and Genentech outside the submitted work; and served as Associate Editor at JAMA Neurology. *A.M.S.* reported grants from the National Institutes of Health, the Bluefield Project to Cure Frontotemporal Dementia, and the Association for Frontotemporal Degeneration; personal fees from Alector, Prevail Therapuetics/Eli Lilly, Passage Bio, Takeda, and the Alzheimer's Drug Discovery Foundation; and other from Datacubed Health (licensing fees) outside the submitted work. *H. J. R.* reported consulting fees from Genentech and Eisai outside the submitted work. Dr Staffaroni reported grants from the National Institutes of Health, the Bluefield Project to Cure Frontotemporal Dementia, and the Association for Frontotemporal Degeneration; personal fees from Alector, Prevail Therapuetics/Eli Lilly, Passage Bio, Takeda, and the Alzheimer's Drug Discovery Foundation; and other from Datacubed Health (licensing fees) outside the submitted work. *C.D.* serves as a consultant to Norvo Nordisk and Eisai Pharmaceuticals. *D.K.J.* has stock holdings in Sage Cerebrovascular Diagnostics, serves as President for Sage Cerebrovascular Diagnostics, and has a patent pending for Serologic assay for silent brain ischemia licensed to Sage Cerebrovascular Diagnostics

658    **Consent Statement**

659         Written informed consent was obtained from all participants or their legal representative. All

660    studies were approved by the Institutional Review Boards of the consortia or ADC in accordance with

661    institutional guidelines and the Helsinki Declaration.

662

663

664     **Key Words**:

665              Memory; Alzheimer's Dementia and Related Disorders; Harmonization; Co-calibration;

666     Neuropsychology; Cognition; Neurodegeneration

667

668     **Legends**

669     **Table 1:** Sample demographics, clinical diagnoses, and scores for UDS-M+ and components

670        Presents mean (standard deviation) or count (percent) of demographic variables, diagnoses, and

671        cognitive scores. Data are presented for the full sample, Development cohort, and DVCID.

672        * The Development and DVCID cohorts are compared via T-tests or Chi2 tests and the group

673        with significantly larger value is indicated with an asterisk. Comparisons were not completed

674        when the sample size of a cohort was small (n<10).

675

676        The "Other" category of clinical diagnoses included unspecified frontotemporal lobar

677        degeneration (n=48), multiple system atrophy (n-1), unspecified primary progressive aphasia

678        (n=14), and those who met no criteria (n=94).

679        Abbreviations: **UDS-M+**: Uniform Data Set Memory plus List-Learning score; **CDR-sb**: Clinical

680        Dementia Rating Scale Sum of Boxes; **CDR-g**: Clinical Dementia Rating Scale Global Score;

681        **MoCA**: Montreal Cognitive Assessment;  **HVLT**: Hopkins Verbal Learning Test; **CVLT-II**:

682        California Verbal Learning Test-Version Two; **CVLT-SF**: California Verbal Learning Test

683        Version Two Short Form; **AVLT**: Rey Auditory Verbal Learning Test; **CERAD**: Consortium to

684        Establish a Registry for Alzheimer's Disease; **PPA:** primary progressive aphasia.

685

686

687    **Figure 1**: UDS-M+ by list-learning task sample within Development cohort

688

689        The histograms display the distribution and mean UDS-M+ score (red dashed line) for

690        subsamples delineated by which list-learning task was administered. All subsamples were derived

691        from the Developmental Cohort. T-tests compared UDS-M+ scores (adjusted for age, Clinical

692        Dementia Rating Scale (CDR) sum of boxes, sex, and education) between list-learning

693        subgroups. All combinations of pairwise comparisons were completed. Cohen's $d$s were

694        calculated for each comparison, and the average Cohen's $d$ is presented above each histogram; for

695        example, -0.085 is the average Cohen's $d$ for pairwise comparisons of the AVLT group and all

696        other list-learning groups.  **Abbreviations**: **AVLT**: Rey-Auditory Verbal Learning Test; **CERAD**:

697        Consortium to Establish a Registry for Alzheimer's Disease; **CVLT-II:** California Verbal

698        Learning Task-II Standard form; **CVLT-SF**: CVLT-II Short form; **HVLT**: Hopkins Verbal

699        Learning Task; **UDS-Only**: Participants who only completed a subset of linking items but were

700        not administered a list-learning task.

701

702 **Figure 2:** Marginal reliability

703 Presents marginal reliability within a sample that combines the Development cohort. The left

704 graph includes all participants. The right graph separates participants based on which list-learning

705 task they completed.  Horizonal dashed lines indicate when the marginal reliability is at 0.8 and

706 0.9. Abbreviations: AVLT: Rey-Auditory Verbal Learning Test; CERAD: Consortium to Establish

707 a Registry for Alzheimer's Disease; CVLT-II: California Verbal Learning Task-II Standard form;

708 CVLT-SF: CVLT Short form; HVLT: Hopkins Verbal Learning Task; UDS-Only: Participants

709 who only completed a subset of anchor items.

710

711 **Figure 3**: Convergent and discriminant validation
712
713 Forest plot presenting standardized regression betas in the Development (circle) and DVCID
714 cohorts (triangle), after controlling for CDR Sum of boxes (CDR Sb), with the exception of the
715 CDR Sb row. Confidence intervals (95%) are shown via the horizontal lines. Variables are color
716 coded by domain.
717 * indicates reverse scoring. Abbreviations: MoCA: Montreal Cognitive Assessment total score;
718 TabCAT: Tablet-based Cognitive Assessment Tool; UDS-EF: Uniform Data Set (v3.0) executive
719 function composite score; DS: Digit Span; Words: Lexical Fluency; MINT= Multilingual Naming
720 Task
721

722    **Figure 4**: Known Group Validity: UDS-M+ by Clinical Dementia Rating Scale
723

724         These graphs illustrate the stepwise decline in UDS-M+ scores that was observed with increasing

725    disease stage (p<.001) in two independent cohorts. Abbreviations:  UDS-M+: Uniform Data Set Version 3

726    Memory plus List-Learning score; CDR: Clinical Dementia Rating Scale

727

728    **Figure 5**: Association between UDS-M+ and plasma phosphorylated tau levels

729

730    Each graph presents the association between UDS-M+ and plasma phosphorylated tau levels

731    (Development cohort: p-tau217; DVCID: p-tau181), after controlling for age, sex, education, and

732    CDR sum of boxes. Graphs on the left present the results for the full sample within Development

733    cohort (A) and DVCID (B). The stacked graphs on the right display the same results, separated by

734    CDR global scores. p-tau levels were log transformed then Z scored. Abbreviations:  UDS-M+:

735    Uniform Data Set Memory plus List-Learning score; CDR: Clinical Dementia Rating Scale

1                                                   Tables and Figures

**Tabel 1: Sample demographics, clinical diagnoses, and scores for UDS-M+ and components**

|  | **Overall** | **Development** | **DVCID** |
|---|---|---|---|
| Sample Size | 7001 | 5943 | 1058 |
| Age (years) | 66.57 (11.19) | 65.34 (11.39) | 73.91 (5.85)* |
| Education (years) | 16.01 (2.75) | 16.04 (2.78) | 15.83 (2.59) |
| Race/Ethnicity |  |  |  |
| White non-Hispanic | 2143 (49.0) | 1787 (50.6)* | 356 (42.5) |
| Hispanic or Latino | 499 (11.4) | 313 (8.9) | 186 (22.2)* |
| Black/African American | 753 (17.2) | 464 (13.1) | 289 (34.5)* |
| American Indian or Alaskan Native | 39 (0.9) | 38 (1.1) | 1 (0.1) |
| Pacific Islander | 300 (8.1) | 299 (10.4) | 1 (0.1) |
| Asian | 865 (19.8) | 864 (24.5) | 1 (0.1) |
| Other | 30 (0.7) | 25 (0.7) | 5 (0.6) |
| Female (%) | 3878 (61.0) | 3297 (59.7) | 581 (69.3)* |
| MoCA Total Score | 22.65 (5.71) | 22.13 (5.99) | 24.91 (3.39)* |
| CDR-g (%) |  |  |  |
| 0 | 3049 (50.2) | 2498 (47.7) | 551 (65.8) |
| 0.5 | 2214 (36.4) | 1927 (36.8) | 287 (34.2) |
| 1 | 693 (11.4) | 693 (13.2) | 0 (0.0) |
| 2 | 106 (1.7) | 106 (2.0) | 0 (0.0) |
| 3 | 15 (0.2) | 15 (0.3) | 0 (0.0) |
| CDR-sb | 1.56 (2.41) | 1.75 (2.53)* | 0.38 (0.68) |
| UDS-M+ | 0.04 (0.95) | -0.03 (0.99) | 0.36 (0.73)* |
| Diagnostic Syndrome | **Overall** | **Development** | **DVCID** |
| Cognitively Unimpaired | 1448 (52.0) | 889 (43.6) | 559 (75.1)* |
| Mild Cognitive Impairment | 435 (15.6) | 258 (12.7)* | 177 (23.8) |
| Early-Onset Alzheimer's Dementia | 182 (6.5) | 182 (8.9) | 0 (0.0) |
| Early-Onset Non-Alzheimer's Dementia | 68 (2.4) | 68 (3.3) | 0 (0.0) |
| Alzheimer's Dementia Syndrome | 211 (7.6) | 211 (10.4) | 0 (0.0) |
| Logopenic Variant PPA | 20 (0.7) | 20 (1.0) | 0 (0.0) |
| Dementia with Lewy Bodies | 9 (0.3) | 9 (0.4) | 0 (0.0) |
| Behavioral Variant FTD | 51 (1.8) | 51 (2.5) | 0 (0.0) |

| | | | |
|---|---|---|---|
| Amyotrophic Lateral Sclerosis | 10 (0.4) | 10 (0.5) | 0 (0.0) |
| Corticobasal Syndrome | 42 (1.5) | 42 (2.1) | 0 (0.0) |
| Non-Fluent Variant PPA | 25 (0.9) | 25 (1.2) | 0 (0.0) |
| Semantic Variant PPA | 38 (1.4) | 38 (1.9) | 0 (0.0) |
| Progressive Supranuclear Palsy | 27 (1.0) | 27 (1.3) | 0 (0.0) |
| Parkinson's Dementia | 6 (0.2) | 6 (0.3) | 0 (0.0) |
| Vascular Dementia | 4 (0.1) | 4 (0.2) | 0 (0.0) |
| Unspecified Dementia | 8 (0.3) | 0 (0.0) | 8 (1.1) |
| Psychiatric | 12 (0.6) | 12 (1.0) | 0 (0.0) |
| Other | 185 (9.0) | 185 (9.0) | 0 (0.0) |

Presents mean (standard deviation) or count (percent) of demographic variables, diagnoses, and cognitive scores. Data are presented for the full sample, Development cohort, and DVCID.
* The Development and DVCID cohorts are compared via T-tests or Chi2 tests and the group with significantly larger value is indicated with an asterisk. Comparisons were not completed when the sample size of a cohort was small (n<10).

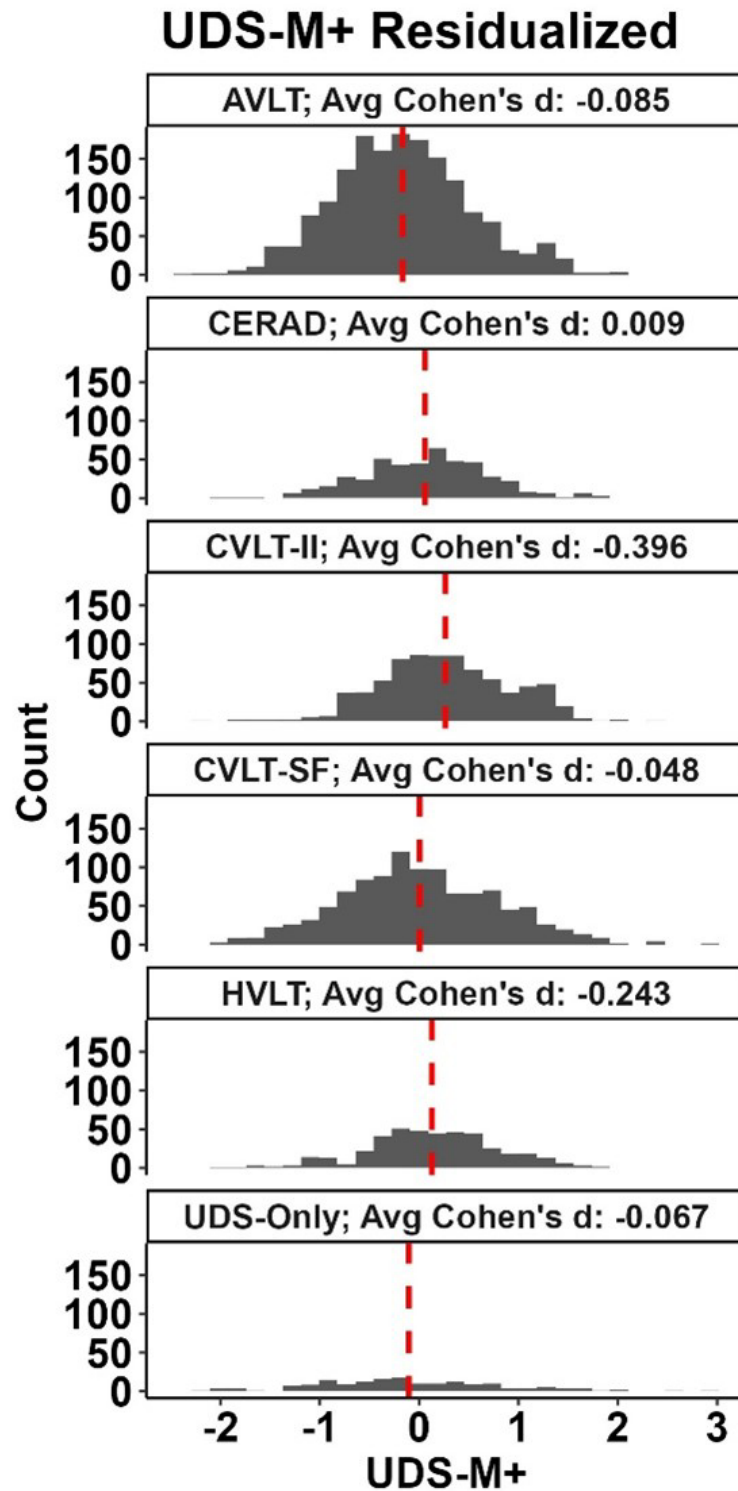The "Other" category of clinical diagnoses included unspecified frontotemporal lobar degeneration (n=48), multiple system atrophy (n-1), unspecified primary progressive aphasia (n=14), and those who met no criteria (n=94).
Abbreviations: **UDS-M+**: Uniform Data Set Memory plus List-Learning score; **CDR-sb**: Clinical Dementia Rating Scale Sum of Boxes; **CDR-g**: Clinical Dementia Rating Scale Global Score; **MoCA**: Montreal Cognitive Assessment;  **HVLT**: Hopkins Verbal Learning Test; **CVLT-II**: California Verbal Learning Test-Version Two; **CVLT-SF**: California Verbal Learning Test Version Two Short Form; **AVLT**: Rey Auditory Verbal Learning Test; **CERAD**: Consortium to Establish a Registry for Alzheimer's Disease; **PPA:** primary progressive aphasia.
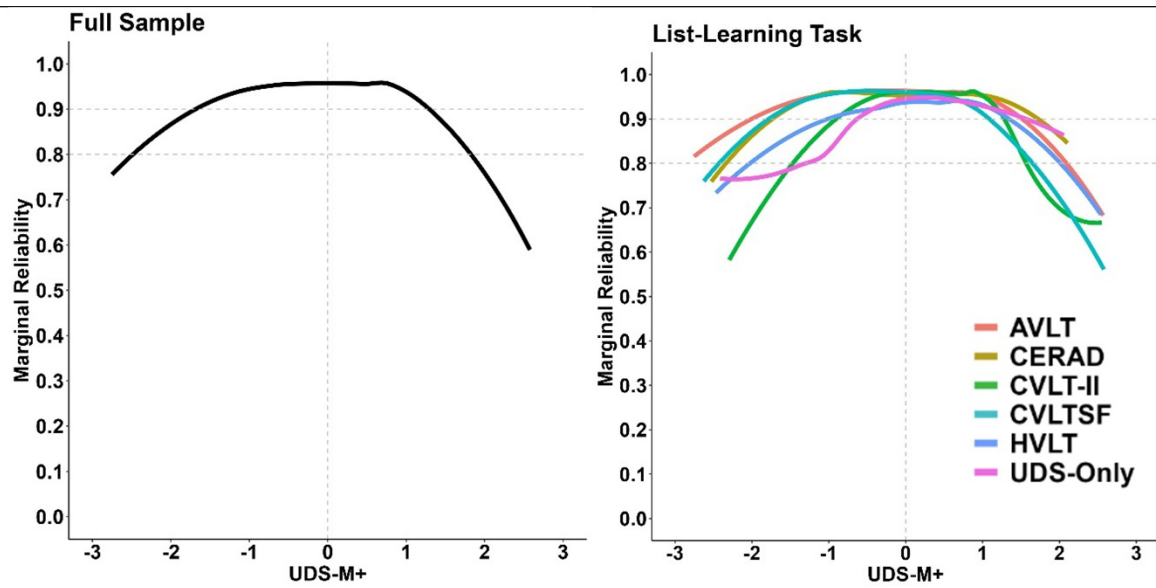
2

3

4    Figure 1

## UDS-M+ Residualized



5

6
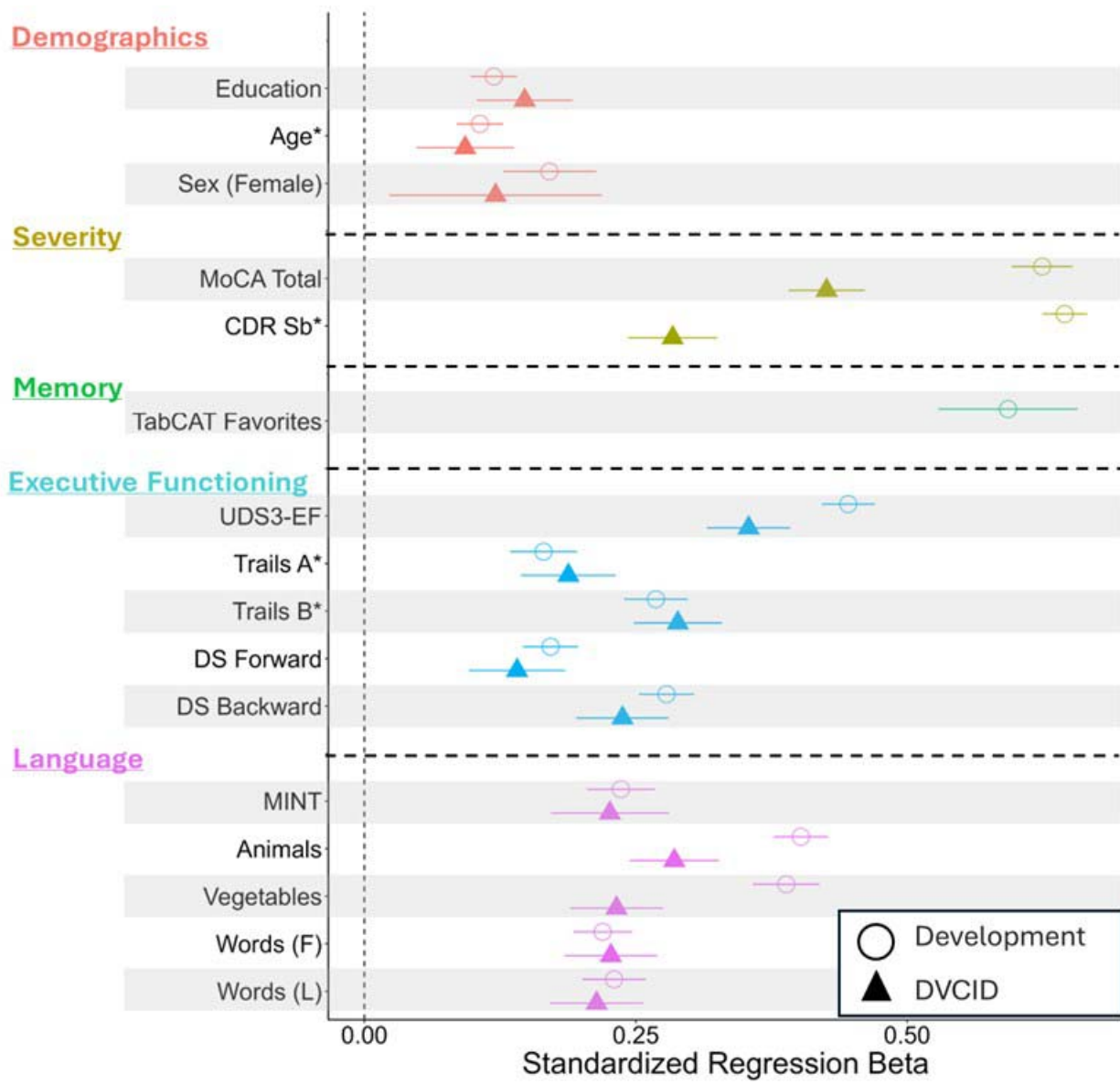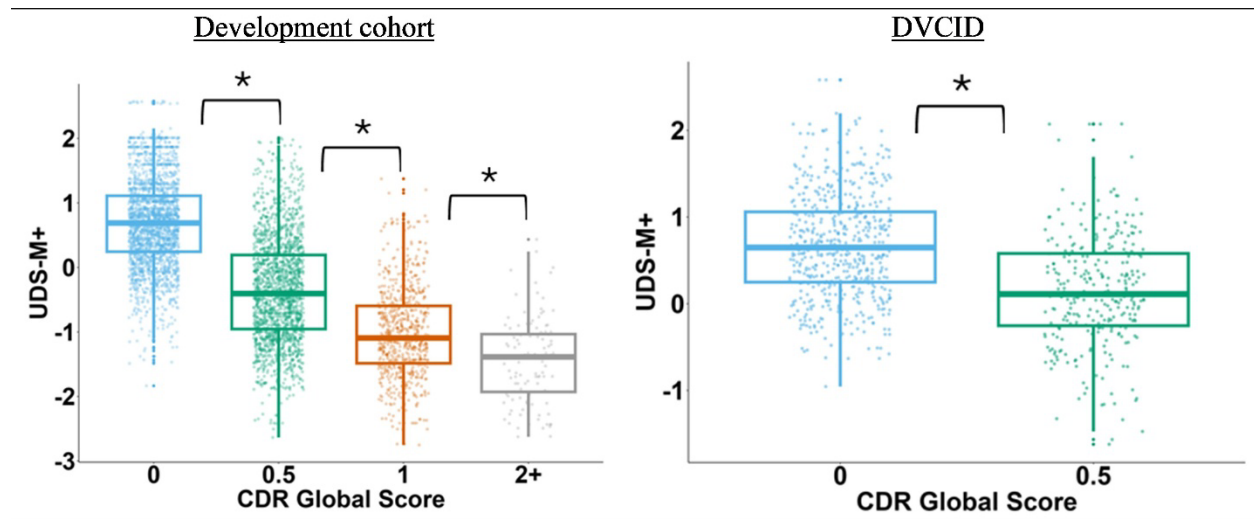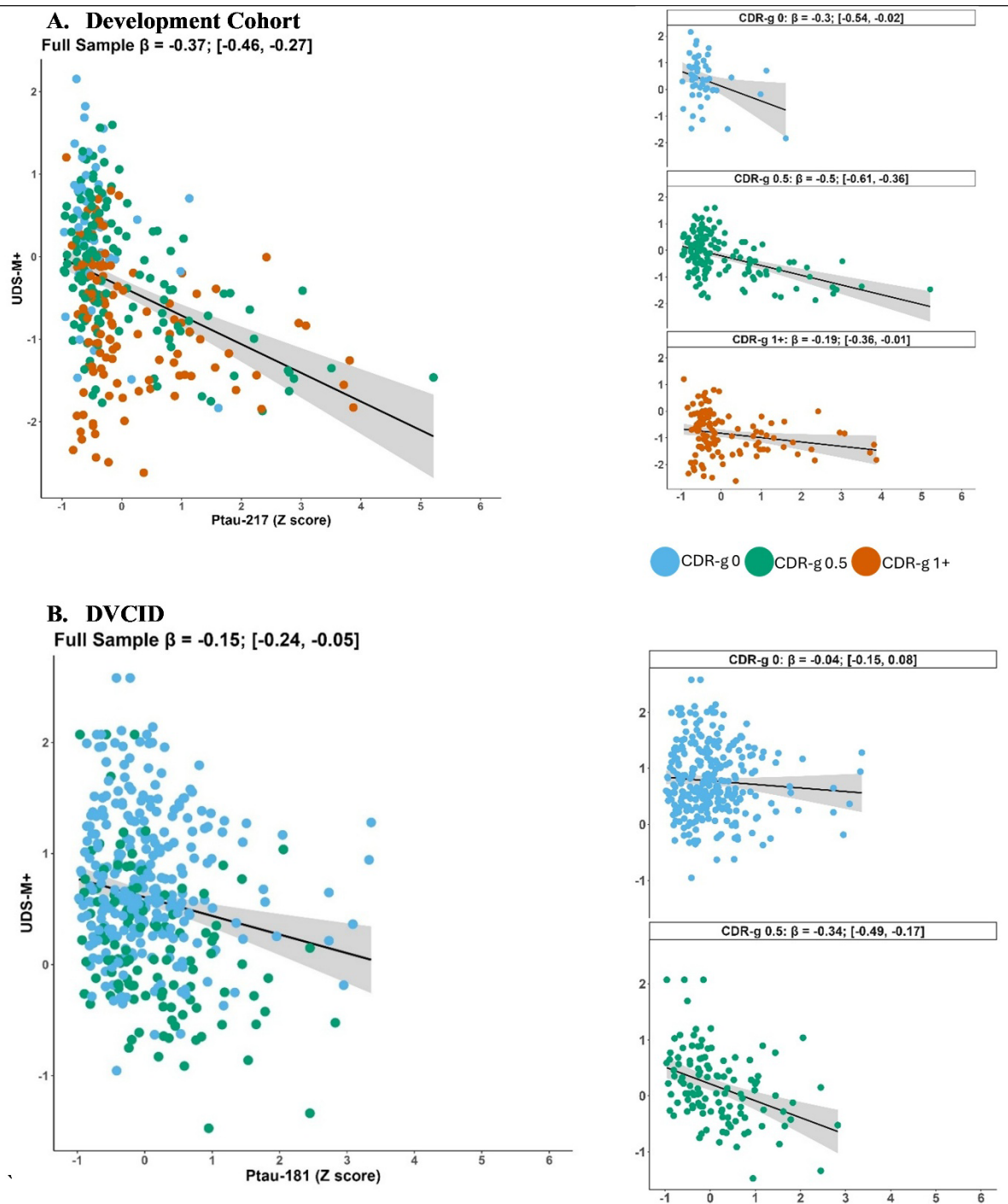
7     Figure 2



8

9

10    Figure 3



11

12

13    Figure 4



14

15

16    Figure 5

19