

1 August 8<sup>th</sup>, 2024

2

3 **Development of a Simplified Smell Test to Identify Patients with Typical**  
4 **Parkinson's as Informed by Multiple Cohorts, Machine Learning and**  
5 **External Validation**

6 Juan Li, PhD<sup>1,2,4,14\*</sup>, Kelsey Grimes, MSc<sup>1,2</sup>, Joseph Saade, MD<sup>1</sup>, Julianna J. Tomlinson,  
7 PhD<sup>1,4,5,14,15</sup>, Tiago A. Mestre, MD, PhD<sup>1,2,4,6,8</sup>, Sebastian Schade, MD<sup>9</sup>, Sandrina Weber, MD<sup>10</sup>,  
8 Mohammed Dakna, PhD<sup>10</sup>, Tamara Wicke, MSc<sup>9</sup>, Elisabeth Lang, BSc<sup>9</sup>, Claudia Trenkwalder,  
9 MD<sup>9</sup>, Natalina Salmaso, PhD<sup>11,14,15</sup>, Andrew Frank, MD<sup>4,12</sup>, Tim Ramsay, PhD<sup>2,3,7</sup>, Douglas  
10 Manuel, MD, MSc, FRCPC<sup>2,6,7</sup>, aSCENT-PD Investigators<sup>14,15</sup>, Brit Mollenhauer, MD<sup>9,10,13,14,15\*</sup>,  
11 Michael G. Schlossmacher, MD, FRCPC<sup>1,4,5,6,8,14,15\*</sup>

12

13 <sup>1</sup> Neuroscience Program, <sup>2</sup> Methodological and Implementation Research Program, and <sup>3</sup> the  
14 Methods Centre, Ottawa Hospital Research Institute, Ottawa, ON., Canada

15 <sup>4</sup> University of Ottawa Brain and Mind Research Institute, Ottawa, ON., Canada

16 <sup>5</sup> Department of Cellular and Molecular Medicine, <sup>6</sup> Department of Medicine, and <sup>7</sup> School of  
17 Epidemiology and Public Health, Faculty of Medicine, University of Ottawa, Ottawa, ON,  
18 Canada

19 <sup>8</sup> Division of Neurology, Department of Medicine, The Ottawa Hospital, Ottawa, ON., Canada

20 <sup>9</sup> Paracelsus-Elena-Klinik, Kassel, Germany

21 <sup>10</sup> Department of Neurology, University Medical Center Goettingen, Germany

22 <sup>11</sup> Department of Neuroscience, Carleton University, Ottawa, ON, Canada

23 <sup>12</sup> Memory Program, Bruyère Research Institute, Ottawa, ON, Canada

24 <sup>13</sup> Deutsches Zentrum für Neurodegenerative Erkrankungen (DZNE), Goettingen, Germany

*Li et al., Simplified Olfaction Testing to Identify Patients with Parkinson's*

25 <sup>14</sup> Aligning Science Across Parkinson's (ASAP) Collaborative Research Network, Chevy Chase,  
26 MD 20815

27 <sup>15</sup> aSCENT-PD Investigators include: Ben Arenkiel, Zhandong Liu, Brit Mollenhauer, Josef  
28 Penninger, Max Rousseaux, Natalina Salmaso, Michael Schlossmacher, Christine Stadelmann,  
29 Julianna Tomlinson, John M. Woulfe.

30 \* Corresponding authors

31 Correspondence to:

32 Dr. Juan Li

33 Ottawa Hospital Research Institute

34 [juli@ohri.ca](mailto:juli@ohri.ca)

35

36 Dr. Brit Mollenhauer

37 Paracelsus-Elena Klinik

38 [brit.mollenhauer@med.uni-goettingen.de](mailto:brit.mollenhauer@med.uni-goettingen.de)

39

40 Dr. Michael Schlossmacher

41 Ottawa Hospital Research Institute

42 [mschlossmacher@toh.ca](mailto:mschlossmacher@toh.ca)

43

44

45 **ABSTRACT**

46 **Background:** Reduced olfaction is a common feature of patients with typical Parkinson disease  
47 (PD). We sought to develop and validate a simplified smell test as a screening tool to help  
48 identify PD patients and explore its differentiation from other forms of parkinsonism.

49 **Methods:** We used the Sniffin' Sticks Identification Test (SST-ID) and the University of  
50 Pennsylvania Smell Identification Test (UPSIT), together with data from three case-control  
51 studies, to compare olfaction in 301 patients with PD or dementia with Lewy bodies (DLB) to 36  
52 subjects with multiple system atrophy (MSA), 32 individuals with progressive supranuclear  
53 palsy (PSP) and 281 neurologically healthy controls. Individual SST-ID and UPSIT scents were  
54 ranked by area under the receiver operating characteristic curve (AUC) values for group  
55 classification, with 10-fold cross-validation. Additional rankings were generated by leveraging  
56 results from eight published studies, collectively including 5,853 unique participants. Lead  
57 combinations were further validated using (semi-)independent datasets. An abbreviated list of  
58 scents was generated based on those shared by SST-ID and UPSIT.

59 **Findings:** We made the following five observations: (i) PD and DLB patients generally had  
60 worse olfaction than healthy controls, as published, with scores for MSA and PSP patients  
61 ranking as intermediate. (ii) SST-ID and UPSIT scents showed distinct discriminative  
62 performances, with the top odorants (licorice, banana, clove, rose, mint, pineapple and cinnamon)  
63 confirmed by external evidence. (iii) A subset of only seven scents demonstrated a similar  
64 performance to that of the complete 16-scent SST-ID and 40-scent UPSIT kits, in both discovery  
65 and validation steps. Seven scents distinguished PD/DLB subjects from healthy controls with an  
66 AUC of 0.87 (95%CI 0.85-0.9) and PD/DLB from PSP/MSA patients with an AUC of 0.73  
67 (95%CI 0.65-0.8) within the three cohorts (n=650). (iv) Increased age was associated with a  
68 decline in olfaction. (v) Males generally scored lower than females, although this finding was not  
69 significant across all cohorts.

70 **Interpretation:** Screening of subjects for typical Parkinson's-associated hyposmia can be  
71 carried out with a simplified scent identification test that relies on as few as seven specific

72 odorants. There, the discrimination of PD/DLB subjects vs. age-matched controls is more  
73 accurate than that of PD/DLB vs. PSP/MSA patients.

74 **Funding:** This work was supported by: Parkinson Research Consortium; uOttawa Brain & Mind  
75 Research Institute; and the Aligning Science Across Parkinson's Collaborative Research  
76 Network.

77 **Key words:** Parkinson disease; parkinsonism; olfaction; hyposmia; smell test; Sniffin' Stick Test;  
78 University of Pennsylvania Smell Identification Test; dementia

79

## 80 **Research in context**

### 81 **Evidence before this study**

82 Chronic hyposmia is a common feature of Parkinson disease (PD) and dementia with Lewy  
83 bodies (DLB), which often precedes motor impairment and cognitive dysfunction by several  
84 years; it is also frequently associated with  $\alpha$ -synuclein aggregate formation in the bulb. The  
85 presence of hyposmia increases an individual's likelihood of having -what has recently been  
86 proposed as- a neuronal synucleinopathy disease, by >24-fold. Despite the strong association of  
87 PD with reduced olfaction, little is understood about it clinically, such as whether it is affected  
88 by sex and age, and whether hyposmia of PD is associated with the same scent identification  
89 difficulty seen in other conditions that present with parkinsonism. Moreover, due to its time-  
90 consuming nature and traditional administration by healthcare workers, extensive olfactory  
91 testing is not routinely performed during neurological assessments in movement disorder clinics.

### 92 **Added value of this study**

93 We analyzed the performance of both the Sniffin' Sticks Test kit and UPSIT battery to  
94 discriminate between healthy controls, patients with PD/DLB and those with MSA or PSP.  
95 Comparison to and juxtaposition with eight other published studies allowed for the generation of  
96 a markedly abbreviated smell identification test that unified both tests, as described. Group

97 classification performance by each scent and its distractors was further analyzed using machine  
98 learning and advanced Item Response Theory methods. Relations between each scent tested, sex  
99 and age were analyzed for the first time. Our findings suggest concrete steps to be implemented  
100 that would allow for simplified, routine olfaction testing in the future.

### 101 **Implications of all the available evidence**

102 Olfaction testing has emerged as an important neurological assessment part when examining  
103 subjects with Parkinson's and those at risk of it. A simple, validated smell test containing fewer  
104 scents than current options could facilitate rapid testing of olfaction in clinic settings and at home,  
105 without supervision by healthcare workers. The usefulness of such a non-invasive test in  
106 population health screening efforts could be further enhanced when coupled to a self-  
107 administered survey that includes questions related to other risk factors associated with PD. As  
108 such, large-scale community screening and applications to routine practice in family doctors'  
109 offices as well as in specialty clinics could be made operationally feasible and cost-effective.

## 110 INTRODUCTION

111 Hyposmia is a common non-motor sign of Parkinson disease (PD) and dementia. The reported  
112 prevalence of olfaction loss in PD ranges from 45% to >90% based on populations selected,  
113 testing methods, and threshold criteria.<sup>[1]</sup> Chronic hyposmia is also described as predictive, with  
114 reduced olfaction preceding PD diagnosis by 4-20 years.<sup>[1]-[3]</sup> Olfactory testing may also help in  
115 the differentiation of parkinsonian syndromes.<sup>[4]</sup> Several screening tools and predictive models  
116 for the incidence of PD have included subjective or objective assessments of olfaction.<sup>[5]-[9]</sup>

117 Two commonly used smell tests for evaluating olfactory function include the University of  
118 Pennsylvania Smell Identification Test (UPSIT),<sup>[10]</sup> a battery of 40 scratch-and-sniff questions  
119 that are self-administered, and the Sniffin' Sticks Test (SST) battery,<sup>[11],[12]</sup> with three subtests  
120 (for Identification; Discrimination; Threshold) comprising 16 scents each, of which the  
121 administration usually involves a trained supervisor. The SST-Identification (SST-ID) and  
122 UPSIT are comparable in that they both assess one's ability to identify a range of scents. Within  
123 SST, SST-ID is more commonly used in PD cohort studies and known to have better diagnostic  
124 performance than the SST-Discrimination subtest.<sup>[13]</sup>

125 Smell test kits were initially developed to assess olfaction in the general population but have  
126 been increasingly used in research settings that study disorders of the brain. Using different  
127 cohorts and methods, some studies have ranked odorants in the UPSIT<sup>[14]-[17]</sup> and SST-ID  
128 kits<sup>[13],[18]-[20]</sup> by their diagnostic performances, and reported that certain subsets of scents  
129 appeared to have equal or better performance than the entire complement of 40- or 16-scents-  
130 based tests. Other studies<sup>[21]-[24]</sup> have examined subsets of UPSIT but without any details on  
131 scent ranking. Moreover, external validation was frequently missing in these analyses, proposed  
132 scent combinations were found to be cohort-specific without agreement across different  
133 studies,<sup>[16],[25]</sup> and the role of distractors (versus the correct scent offered) in such multi-choice  
134 settings were understudied. Furthermore, analyses of UPSIT and SST-ID kits were always  
135 conducted separately despite the similarities between the two tests. Finally, olfaction scores in  
136 patients with other, atypical forms of parkinsonism were not assessed in PD-centric studies.

137 In the current work, we aimed to assess olfaction performances in commonly encountered forms  
138 of parkinsonism, to assess the key features of both UPSIT and SST-ID odorants, to explain any  
139 observed differences between them, and to develop a simplified smell test to unify both kits  
140 using proper internal and external validation steps. To this end, eight published scent  
141 rankings,<sup>[13]-[20]</sup> which collectively examined 5,853 participants, were incorporated into our study  
142 to make the proposed abbreviated test generalizable and to avoid overfitting. We also added Item  
143 Response Theory-based analyses when examining the behaviour of participants' responses to  
144 multiple choices provided for each scent; this, to enhance sensitivity and specificity of future test  
145 versions. Lastly, we analyzed the effects of age and sex on performance in scent identification.

146

## 147 **METHODS**

148 The study was conducted in adherence with STARD<sup>[26]</sup> guideline, see **Supplemental Table 1**.

### 149 *Source of data and participants*

150 We used de-identified data from three observational, retrospective, case-control studies: the “*De*  
151 *Novo* Parkinson disease study” (DeNoPa);<sup>[27]</sup> Ottawa (PREDIGT) Trial; and Prognostic  
152 Biomarkers in Parkinson's Disease Study (PROBE).<sup>[28]</sup> Detailed description and  
153 inclusion/exclusion criteria of these three cohorts are provided in **Supplemental Methods**; their  
154 demographic and diagnostic characteristics are summarized in **Table 1**. Data of the cross-  
155 sectional Ottawa Trial study, baseline data of the longitudinal PROBE study, and three visits of  
156 the longitudinal DeNoPa study (baseline; 48-month; and 72-month follow-up visits) were used.  
157 Patients with PD or dementia with Lewy bodies (DLB) (DeNoPa: n=129; Ottawa Trial: n=70;  
158 PROBE: n=102), subjects with multiple system atrophy (MSA) or progressive supranuclear  
159 palsy (PSP) (DeNoPa: n= 9; Ottawa Trial: n=6; PROBE: n=53), and neurologically healthy  
160 controls (HC) (DeNoPa: n=109; Ottawa Trial: n=118; PROBE: n=54) were included. Diagnostic  
161 criteria were applied according to international guidelines;<sup>[27]-[30]</sup> diagnoses were revised, as  
162 necessary, during longitudinal follow up visits or by independent raters who were blind to

163 olfaction test results. Most study participants with PD in the three trials were classified as Hoehn  
164 and Yahr stage II-III. No participant overlap existed between the three studies.

165 Analyses of these deidentified cohort data were approved by Investigational Review Boards at  
166 Paracelsus-Elena-Klinik (Kassel, Germany) in Frankfurt, Hesse (FF 89/2008), at The Ottawa  
167 Hospital (Ottawa, Canada; 20180010-01H) as well as all PROBE Study-affiliated sites, with the  
168 participants' consent.

### 169 *Study assessments and statistical analyses*

170 *Sniffin' Sticks test (SST)*: Used in DeNoPa, the SST comprises a supervised test for one's sense  
171 of smell using pen-like odor dispensing devices, as administered in a clinic setting.<sup>[11],[12]</sup> The test  
172 has three subtests: Identification (SST-ID), Discrimination (SST-DS), and Threshold (SST-TH),  
173 each with 16 odorants. In SST-ID, subjects are presented a stick and choose the scent from four  
174 options (one correct, three distractors). SST-DS is performed using triplets of odorants that are of  
175 similar intensity and hedonic tone, where subjects are required to identify which stick of the  
176 triplet has a different scent from the other two. SST-TH is performed using triplets of sticks  
177 where only one is filled with odorant at a certain dilution whereas the other two are filled with  
178 odor-free solvent. SST-TH determines at what dilution subjects can consistently identify the  
179 odorant-filled stick. The entire SST (in German) was completed by all DeNoPa participants at  
180 their baseline visit, and the SST-ID subtest was re-administered at 48-month and 72-month  
181 follow-up visits.

182 *The University of Pennsylvania Smell Identification Test (UPSIT)*: UPSIT was used in the  
183 Ottawa Trial and PROBE; this self-administered kit contains 40 scratch-and-sniff questions; each  
184 question has one correct answer and three incorrect distractors.<sup>[10]</sup>

185 *Data preparation*: Observations that had no valid response to SST-ID or UPSIT were removed  
186 from the analysis. Observations with incomplete responses were imputed with 0s, indicating  
187 incorrect responses. A dichotomous response-based transformation (0 = incorrect response; 1 =  
188 correct response) was used to calculate the sum scores and assess discriminative performances



189 for each scent. The exact indices of chosen options were used for Item Response Theory analysis  
190 (see below).

191 *Demographic and diagnostic characteristics:* Demographic and diagnostic characteristics of the  
192 study cohort were summarised using n (%) or median (the interquartile range (IQR)). The  
193 reported p-values represented the significance from corresponding Fisher's exact test or Kruskal-  
194 Wallis rank sum test, with false discovery rate correction for multiple testing; p-values smaller  
195 than 0.05 were considered significant.

196 *Comparing different smell tests:* Score distributions of corresponding smell tests in each subject  
197 group were illustrated using Cummings estimation plots.<sup>[31]</sup> The raw UPSIT and SST-ID scores  
198 were also normalized into percentiles based on age and sex, and hyposmia was defined by SST-  
199 ID percentile  $\leq 10\%$ <sup>[32]</sup> or UPSIT percentile  $\leq 15\%$ .<sup>[10]</sup> Discrimination performances of these  
200 subtests were compared using area under the receiver operating characteristic (ROC) curve  
201 (AUC) values with bootstrap estimated 95% confidence interval (CI),<sup>[33]</sup> in order to distinguish  
202 diagnostic groups. **Table 2** also reports optimal thresholds and their associated sensitivity and  
203 specificity that correspond to the maximum Youden indices.<sup>[34]</sup>

204 *Machine learning workflow of developing and validating the abbreviated smell test:* **Figure 1**  
205 illustrates the machine learning workflow of ranking UPSIT/SST-ID scents, developing and  
206 validating respective simplified tests and the unified abbreviated test. Data of the Ottawa Trial  
207 and baseline data of DeNoPa were used as discovery cohorts, and baseline data of PROBE and  
208 follow-up data of DeNoPa were used for (semi-)independent validation.

209 1. *Ranking the individual scents:* Using the corresponding discovery cohorts with 10-fold  
210 cross validation (see **Supplemental Methods**), individual scents in SST-ID and UPSIT  
211 were ranked separately based on their AUC values in differentiating PD/DLB patients  
212 from healthy controls. To control over-fitting, the SST-ID and UPSIT scent rankings  
213 from our study was compared with eight external rankings,<sup>[13]-[20]</sup> four for each test, and  
214 two final lists were generated by averaging internal and external rankings. Eleven scents  
215 are shared by both smell tests, and an additional “Shared” ranking was constructed using  
216 their respective positions in the averaged SST-ID and UPSIT rankings.

217 2. *Developing and validating the best-performing, simplified tests:* For each scent ranking,  
218 beginning with the highest-ranked odorant, subsets were constructed by adding one scent  
219 at a time in descending ranking order. A total of 95 and 220 distinct SST-ID and UPSIT  
220 subsets with various numbers of scents were compared, using their AUC values in  
221 distinguishing PD/DLB from HC, to develop the best-performing simplified tests,  
222 including one that unified both smell tests. The resulted abbreviated smell tests were also  
223 validated using (semi-)independent datasets.

224 *Exploring observed differences in scent performance:* The percentage of correct scent  
225 identification within each subject group was calculated. These percentages were further  
226 compared to examine the relationship of scent identification versus sex and age (using spline  
227 smoothing, see **Supplemental Methods**). Furthermore, Item Characteristic Curves (ICCs) for  
228 each scent within PD/DLB and HC groups from baseline DeNoPa and Ottawa Trial visits were  
229 used to analyze scent performance and the influence of distractors (see **Supplemental Methods**  
230 for more details).

231 Statistical analyses were performed using 'R' (version 4.3.1). Cummings estimation plots were  
232 generated using 'dabestr',<sup>[36]</sup> and ICC curves were generated using 'TestGardener';<sup>[37]</sup> all other  
233 plots were generated using 'ggplot2'.<sup>[38]</sup> The package 'pROC'<sup>[33]</sup> was used for ROC and AUC,  
234 and 'fda' was used for spline smoothing.<sup>[39]</sup>

235

## 236 **RESULTS**

### 237 **Comparison of different smell tests for classifying typical Parkinson disease**

238 We first assessed the performance of UPSIT and SST-ID kits in each diagnostic group. As  
239 expected, across all three cohorts, PD and DLB patients generally had lower scores (*i.e.*, worse  
240 olfaction) than healthy controls (HC), whereas scores of MSA and PSP patients were  
241 intermediate; these are shown by score distributions in **Figure 2(a)-(d)**, and median  
242 scores/percentiles and percentage of hyposmia in **Table 1**. There was no significant difference in  
243 olfaction performance between MSA and PSP patients (**Figure 2(b), (d)**). UPSIT and SST-ID

244 kits had comparable performance in distinguishing PD/DLB patients from HC subjects (**Figure**  
245 **2(e)**, left) with AUC values in the three cohorts ranging between 0.89-0.93. Further, both tests  
246 showed reduced performance when distinguishing PD/DLB patients from MSA/PSP patients  
247 (**Figure 2(e)**, right), but with a larger variation (0.69 to 0.92) across the three cohorts, likely due  
248 to the small sample sizes of MSA/PSP groups in the Ottawa Trial and DeNoPa cohorts. When  
249 compared to the other two SST subtests, SST-ID was found to be the most discriminative one in  
250 distinguishing PD/DLB from both HC and from MSA/PSP (**Supplemental Figure 1**), as  
251 expected.<sup>[13]</sup> For cohort-specific thresholds and their corresponding sensitivity and specificity,  
252 see **Table 2**.

253

#### 254 **Performances of individual scents differ in discriminating PD/DLB from HC**

255 We next sought to determine whether subsets of scents were more informative to discriminate  
256 between PD/DLB and controls compared to the complete 16-scent (SST-ID) or 40-scent (UPSIT)  
257 tests. **Figure 3(a)** shows the distribution of AUC values for each SST-ID scent across 10 folds  
258 using baseline data from the DeNoPa cohort. Clusters of scents identified included banana and  
259 mint as the two most discriminative scents (individual AUC values,  $\geq 0.725$ ), followed by anise,  
260 coffee, licorice, fish and rose in the second-most discriminative cluster. Compared with SST-ID  
261 scents, clustering was less obvious for the UPSIT scents (**Figure 3(c)**), whose AUC values  
262 ranged between 0.5 to 0.77. For the Ottawa Trial cohort, the top-7 UPSIT scents for identifying  
263 PD/DLB patients included rose, wintergreen, root beer, licorice, dill pickle, mint, and grass.

264 The observed differences in each scent's discriminative performance were further examined by  
265 visualizing the percentages of correct scent identification within each diagnostic group (**Figure**  
266 **3(b)**, (d)) and by the percentage differences between HC and PD/DLB groups (**Supplemental**  
267 **Figure 2**). Regardless of the study cohort and smell test used, PD/DLB patients showed  
268 significantly lower percentages of correctly identifying each scent than control subjects. Scents  
269 that were easy to identify in the HC group but difficult for the PD/DLB group (*i.e.*, larger  
270 percentage differences as in **Supplemental Figure 2**) had larger single-scent AUC values.

271 Scents had poorer discriminative performances when both HC and PD/DLB groups found them  
272 easy (e.g. SST-ID: orange, UPSIT: leather) or difficult (e.g. SST-ID: apple, UPSIT: lemon).

273 Therefore, two rankings for scents from the SST-ID and UPSIT kits were constructed. To  
274 overcome the inherent risk of overfitting due to the modestly sized cohorts used, external  
275 evidence was introduced. **Figure 4** compared the scent rankings from this study with previously  
276 published studies, and two “Average” rankings were derived. For the SST-ID kit, different  
277 studies -despite their differences in cohort design and methods applied (**Supplemental Table 1**)-  
278 showed consensus in that anise, licorice, mint, banana, coffee, fish and rose were the most  
279 discriminative scents in distinguishing PD/DLB subjects from HCs (**Figure 4(a)**). For the UPSIT  
280 battery, however, its related studies showed less agreement on their scent rankings (**Figure 4(b)**),  
281 which could be partially explained by results shown in **Figure 3(c)**; there, many UPSIT scents  
282 showing similar performances and revealed fewer clusters than did SST-ID-based scents.  
283 Nonetheless, the top-7 UPSIT scents in the final “Average” ranking were coconut, clove,  
284 wintergreen, banana, licorice, grass and cherry. With respect to a possibly common, simplified  
285 list, we noted that there are eleven scents shared by SST-ID and UPSIT, and thus, an additional  
286 “Shared” ranking was generated to construct a unified, abbreviated smell test (**Table 3**).

287

### 288 **Item Characteristic Curves further reveal details for each scent and the influence of** 289 **distractors**

290 The findings above revealed subsets of scents that were relatively discriminative (from a PD  
291 perspective), which could suggest a disease-specific and/or scent processing-related change that  
292 is linked to *participant performance*. However, for multiple-choice tests like SST-ID and UPSIT  
293 kits, the selection of distractors paired with each scent could also influence *odorant performance*.  
294 Item Characteristic Curves (ICCs) can help address this, especially when the scent is shared by  
295 different tests.

296 In the current context, mint and licorice were two well-performing scents, and their ICCs  
297 (**Figure 5(1)-(4)**) showed similar characteristics in that HCs generally correctly identified them,

298 while PD/DLB patients had more difficulty in choosing the right option. However, there were  
299 also some noteworthy differences: when scoring on the scent for 'mint', PD/DLB patients could  
300 rule out 'chive' and 'onion' in SST-ID and 'fruit punch' in UPSIT, indicating that they detected  
301 some scent, but it was not declarative enough to choose 'mint'. However, for 'licorice',  
302 particularly in the UPSIT kit, there was strong evidence of random guessing whereby patients  
303 couldn't detect any scent to help favor or eliminate an option. Here, ICCs of scoring by HCs also  
304 eliminated the possibility of the corresponding pen (SST-ID) or encapsulated sticker (UPSIT)  
305 being defective.

306 The scent for 'banana' ranked 1/16 in DeNoPa but only 34/40 in the Ottawa Trial (**Figure 3**);  
307 however, these inconsistent performances were not due to differences in distractors. The option  
308 'cherry' distracted many PD/DLB patients and HCs in the Ottawa Trial, but not in the DeNoPa  
309 study (**Figure 5(5)-(6)**). Here, cohort- or odorant (*e.g.*, its concentration or composition for the  
310 artificial scent)-related differences might be more plausible explanations.

311 'Orange' and 'lemon' were both ranked low in the two tests but for different reasons (**Figure**  
312 **5(7)-(10)**): 'orange' in SST-ID was too easy, even for hyposmic PD/DLB patients. 'Orange' in  
313 UPSIT, however, had different distractors that were active within PD/DLB ('bubble gum') and  
314 HC ('turpentine'). For 'lemon', the distractors of 'grapefruit' in SST-ID and 'motor oil' in  
315 UPSIT confused both patients and healthy persons. Such ICC results within the normosmic  
316 control group might be evidence of a flawed odorant choice or an explanation that is rooted in  
317 chemical manufacturing of the scent. In **Supplemental Figures 3-5** we listed the ICCs for all  
318 other scents.

319

## 320 **Development and validation of abbreviated smell tests**

321 Based on scent rankings in **Figure 4** and **Table 3**, 95 SST-ID subsets and 220 UPSIT subsets  
322 were compared by their AUC values in distinguishing PD/DLB patients from HCs. **Figure 6**  
323 shows the results within each internal and external validation set. When using an increasing  
324 number of highly rank-ordered scents, we observed that the corresponding AUC values for

325 odorant subsets increased steeply for the first four; surprisingly, any improvement in  
326 performance thereafter was marginal. Compared with other published rankings, the “Average”  
327 rankings as well as their subsets appeared to be more discriminative with robust performance in  
328 all the validation sets. Considering a potential trade-off between subset performance and number  
329 of scents administered, the SST-ID subset with seven scents and UPSIT subset with ten scents,  
330 based on their corresponding “Average” rankings, emerged as the best-performing test batteries.

331 With respect to potentially developing a unified smell test that could be applied to large study  
332 populations (those examined here were based in North America and Europe of predominantly  
333 White ethnicity), the subset of seven scents from the “Shared” ranking (**Table 3**) with the highest  
334 performance in all validation sets comprised licorice, banana, clove, rose, mint, pineapple and  
335 cinnamon. Mean (standard deviation) scores and AUC values of this subset within each cohort  
336 are reported in **Table 4**. Note, due to inherent cohort differences, the AUC value in the  
337 diagnostic classification of PD/DLB versus HC when combining all three cohorts was 0.87  
338 (95%CI 0.85-0.9), *i.e.*, slightly lower than those of each separate cohort. In the combined cohort,  
339 the cut-off value for distinguishing PD/DLB versus HC that corresponded to the maximum  
340 Youden index was 4.5 (score range 0-7), the resulted sensitivity and specificity were 0.76 and  
341 0.85, respectively.

342

### 343 **Performances of scents in discriminating PD/DLB from MSA/PSP**

344 The same workflow was applied to the PROBE cohort to generate a subset of scents specialized  
345 in distinguishing between PD/DLB versus MSA/PSP patients. There, a subset with ten scents  
346 (clove, dill pickle, cinnamon, soap, rose, pizza, root beer, turpentine, gasoline, licorice) achieved  
347 an AUC value 0.78 in the validation set within PROBE, a promising improvement when  
348 compared with the entire 40-scent UPSIT test (AUC = 0.68; **Supplemental Figure 6**).  
349 Independent validation is needed to retest the usefulness of this subset given the small number of  
350 MSA/PSP subjects enrolled in the other two cohorts studied herein.

351



## 352 **Assessment of age and sex on scent identification**

353 We also investigated the influence of age and sex on scent identification. **Supplemental Table 2**  
354 shows the coefficients of the linear regression for the relationship between smell test score with  
355 age, sex, and diagnostic groups within each cohort. Not surprisingly, progression in age  
356 significantly lowered olfaction across all groups, and males generally had a worse sense of smell  
357 than their female counterparts, although the latter was not significant across the three cohorts.

358 When focusing just on the eleven scents shared by both tests, relationships between scent  
359 identification and sex were further evaluated by comparing the percentages of correct scent  
360 identification across groups (**Figure 7(a)**). In line with the regression results, females showed  
361 higher percentages of correct identification than males for most of the scents, except for  
362 cinnamon, turpentine, and leather. Next, we compared the probability of identifying each scent  
363 correctly across ages between the PD/DLB and HC groups (**Figure 7(b)**). As expected, older  
364 participants showed decreasing percentages for identifying specific scents correctly. The fitted  
365 lines for PD/DLB and HC groups were usually in the same direction and of similar slopes, with  
366 some exceptions, but these were not consistent across all three cohorts.

367

## 368 **DISCUSSION**

369 To our knowledge, this is the most comprehensive study to date describing olfactory dysfunction  
370 in late-onset, typical PD and two other neurological disorders presenting with parkinsonism  
371 using both the multimodal SST battery and UPSIT kit. The principal take home message from  
372 this study is that when probing for hyposmia in PD, the following points matter: PD/DLB  
373 patients had worse olfaction than healthy subjects, and scores of MSA/PSP patients were  
374 intermediate; there was no observed difference in olfaction between MSA and PSP patients;  
375 scent identification testing is sufficient, and threshold as well as discrimination testing could be  
376 omitted when screening populations for PD using the SST kit; fewer scents can reduce  
377 examination time and test taking fatigue without sacrificing diagnostic accuracy; the selection of  
378 specific scents should be informed by their discriminative performance in specific group

379 classification efforts; random guessing could lower diagnostic accuracy; and from a test design  
380 perspective, choices provided as distractors influence scent performance. Importantly, we found  
381 that a simplified smell test, with specific scents, is sufficient to identify PD/DLB-linked  
382 hyposmia. Such a test, which is now being piloted by us, holds the potential to facilitate olfactory  
383 testing in the clinical setting, used for at-home testing and population-based screening methods.

384 In developing and validating a simplified smell test for this purpose, we used a machine learning  
385 approach and found that only seven scents (licorice, banana, clove, rose, mint, pineapple, and  
386 cinnamon) were required to approximate the diagnostic performance of administering the 16-  
387 scent SST-ID or 40-scent UPSIT batteries and the value of adding more scents was negligible.  
388 Moreover, a test kit for rapid screening with nearly similar performance could be constructed  
389 using as few as four scents.

390 We demonstrated the impact distractors have on detecting specific scents using IRT analysis. We  
391 uncovered uncertainty in eliciting a choice for some scents, even for HCs with intact olfaction.  
392 This could be explained by the difficulty of biological scent discrimination or the to-be-improved  
393 selection of artificial odorants. By extension, our analyses revealed the opportunity to remove ill-  
394 performing scents, *e.g.*, orange and lemon, from currently used kits.

395 Unexpectedly, we found a high level of guessing among PD patients for licorice, indicating  
396 patients' difficulty in detecting this scent. SST-ID and UPSIT batteries are multiple choice-based  
397 tests, in which participants are instructed to always choose even when they cannot smell  
398 anything; such random guessing will introduce errors into data sets. Advanced IRT methods can  
399 treat missing responses as an additional option; administrators of tests would then prefer the  
400 participants to leave any uncertain questions unanswered rather than forcing a guess. However,  
401 for the future administration of olfaction tests, or for designing a new one, we would suggest  
402 adding an extra choice, such as "I cannot identify the scent" to reduce random guessing. Based  
403 on our experience in administering smell tests, the extra option would also help improve  
404 participant experience and eliminate frustration, especially for patients with severe hyposmia.

405 An easy-to-administer, inexpensive, sensitive and non-invasive smell test (with 4-7 scents) could  
406 have important clinical usefulness, particularly when coupled with a short, self-administered



407 questionnaire capturing demographic information and known risk factors of developing PD.<sup>[9]</sup>  
408 Such a questionnaire could also explore other factors leading to hyposmia unrelated to  
409 neurodegeneration, *e.g.*, previous nasal injuries, microbial infections, seasonal allergies, and  
410 chronic exposure to air pollution, to augment specificity for PD. Upon validation, such a kit  
411 could be used as the initial step of large-scale community screening, or in routine clinic practice  
412 of a movement disorders-oriented clinic, or for early detection within a family medicine office.  
413 When it comes to screening efforts for PD, more invasive and expensive tests, *e.g.*, the  $\alpha$ -  
414 synuclein seeding amplification assay from cerebrospinal fluid (CSF), or the administration of a  
415 dopamine transporter scan, could be administered as an additional step to increase screening  
416 accuracy further, such as when considering enrolling PD subjects into specific, disease-  
417 modifying clinical trials.<sup>[40]-[44]</sup>

418 Despite the findings regarding scent ranking and subset analyses, it remains unclear whether a  
419 specific PD olfaction deficit exists, rather than a global reduction in olfaction, and what the  
420 underlying mechanisms would be. We and others recently found that olfactory deficits were  
421 significantly associated with positivity on the  $\alpha$ -synuclein seeding amplification assay in CSF  
422 samples, suggesting that patients may have an underlying disease linked to the dysregulation of  
423 *SNCA* expression and/or protein processing.<sup>[45],[46]</sup>

424 Mechanistically, it remains unknown as to how chronic hyposmia arises in PD/DLB (and REM  
425 Sleep Behaviour Disorder) as well as some MSA/PSP subjects, at what age it begins, the cause  
426 and its underlying circuit-based and molecular mechanisms, and whether olfactory deficits are  
427 shared for specific scents among persons with typical PD versus those with dementia. Large  
428 scale population screening, including with a simplified testing battery derived from SST-ID and  
429 UPSIT kits, could begin to answer some of these questions.

430

## 431 **STRENGTH**

432 In this work, olfaction performance using SST and UPSIT batteries was studied in detail. Both  
433 internal and external validation efforts were used to test/retest performance and avoid overfitting.

434 Scent ranking derived from this study was also compared with eight previously published studies.  
435 By incorporating results from external studies, the unified abbreviated olfaction test kit of using  
436 just seven scents emerged as very generalizable. The observed performance differences for each  
437 scent in group classification was explained using complementary techniques and further studied  
438 considering sex and age. Our study also explored scent identification performance at the level of  
439 choices provided, highlighting the importance of distractors, which allows for future  
440 improvement in the design of test kits.

441

## 442 **LIMITATIONS AND FUTURE WORK**

443 Performance of the unified abbreviated smell test was estimated by simulation: responses of the  
444 related scents were partitioned from the SST-ID and UPSIT data sets, and the original distractors  
445 from SST-ID and UPSIT kits were used. Real performance of the abbreviated test, when  
446 manufactured as a stand-alone product with distractors determined for each scent, should be  
447 assessed in the same cohorts, in new studies, and in routinely run clinic settings. Further, more  
448 data are needed to validate the scent ranking and the associated subset developed here for  
449 distinguishing PD/DLB versus MSA/PSP patients.

450 The three cohorts used are highly homogeneous cohorts with most participants being white.  
451 Although scent rankings and the selection of a simplified smell test have been rigorously  
452 developed and validated with external information incorporated, future calibration and cultural  
453 adaption efforts will be necessary when applying them to other populations, especially those that  
454 differ from Western Europeans.

455 Case-control studies have an inherent potential for selection bias in their recruitment. Especially  
456 because of the age- and sex-matched design, age- and sex-effects were likely underestimated.  
457 Population studies, such as in the initial community screening effort undertaken by PARS  
458 planners and the 'PPMI hyposmia' effort, could provide complementary data sets; however, they  
459 too have potential setbacks: as the majority of participants will have a normal sense of smell, the  
460 score distribution could be highly skewed; a low percentage of PD patients will make the

461 resulting dataset imbalanced; when smell test data are reduced to a single sum score, sub-  
462 analyses will be difficult to complete, which will limit interrogations of data between established  
463 cohorts.

464 For screening purpose, a one-time administered smell test may not be informative enough to  
465 assess a subject's sense of smell completely, because other factors, such as temporary olfaction  
466 reduction/loss due to infection, seasonal allergies, occupational exposure and/or drinking, eating,  
467 smoking before taking the test, could skew results. Retesting at appropriate time intervals will be  
468 required for even higher accuracy.

469 **ACKNOWLEDGEMENTS**

470 The authors acknowledge the commitment of study participants in the three cohorts and are  
471 grateful to all clinical research coordinators at all the study sites. Responsible supervisors for the:

472 DeNoPa Study: Paracelsus-Elena-Klinik: Brit Mollenhauer.

473 Ottawa Trial Study: The Ottawa Hospital: Michael Schlossmacher, Élisabeth Bruyère Hospital:  
474 Andrew Frank.

475 PROBE Study: PROBE Steering Committee: Voyager Therapeutics: Bernard Ravina; Brigham  
476 and Women's Hospital: Clemens Scherzer, University of Ottawa: Michael Schlossmacher, Avid  
477 Radiopharmaceuticals: Andrew Siderowf, University of Rochester: David Oakes, Arthur Watts;  
478 Institute for Neurodegenerative Disorders: Kenneth Marek; Georgetown University: Ira  
479 Shoulson.

480 We thank Nathalie Lengacher for help in graphic design; Nadine Mauri, Nancy MacDonald, and  
481 Yoobin Lee for help in data management/input. This work was supported by funding from  
482 Parkinson Canada (to T.M., D. M., M.G.S; 2018; to J.L.; 2019-2021), Michael J. Fox Foundation  
483 for Parkinson's Research (to T.M., D. M., M.G.S), Department of Medicine (T.M., T.R., D.M.,  
484 M.G.S.), The Ottawa Hospital Foundation (Borealis Foundation to J.L.) and the Uttra & Subhash  
485 Bhargava Family (M.G.S.), the Paracelsus-Elena-Klinik Kassel, Parkinson Fonds Deutschland,  
486 and the Deutsche Parkinson Vereinigung (B.M.; C.T.). The study was also funded by the joint  
487 efforts of The Michael J. Fox Foundation for Parkinson's Research (MJFF) and the Aligning  
488 Science Across Parkinson's (ASAP) initiative. MJFF administers the grant [Grant ID: ASAP-  
489 020625] on behalf of ASAP and itself.

490 The funders had no role in the design and execution of the study; the collection, management,  
491 analysis, and interpretation of the data; the preparation, review, or approval of the manuscript;  
492 and the decision to submit the manuscript for publication. We are grateful for the ongoing  
493 support and feedback from people with lived experiences, such as through the board of the  
494 Parkinson's Research Consortium Ottawa and members of Partners Investing in Parkinson's  
495 Research, and to Drs. P. Wells and D. Lewis for their ongoing encouragement.

496

497 **AUTHOR CONTRIBUTIONS**

498 JL and MGS contributed to the concept and design of the study; JL, KG, and MGS contributed to  
499 the acquisition of data. JL decided on the statistical methods used in this study. JL did data  
500 cleaning, data analysis, figures and tables. JL, JTT, and MGS contributed to data interpretation.  
501 SS, SW, MD, TW, EL, CT, and BM contributed to the data collection and verification of  
502 DeNoPa; KG, JS, JL, JTT, AF, and MGS contributed to the data collection and verification of  
503 Ottawa Trial. JL and MGS wrote the first draft of the manuscript. JTT, BM, NS, TAM, TR, and  
504 DM contributed to the drafting of the manuscript and revising it critically, and all authors  
505 approved the submission of its current version.

506

507 **CONFLICT OF INTEREST STATEMENT**

508 In 2024, MGS co-founded NeuroScent Inc to develop a home-based testing platform for  
509 hyposmia.

510

511 **CODE AVAILABILITY**

512 The code for data analyses and figures is publicly accessible in GitHub (link to be updated).

513

514 **DATA AVAILABILITY**

515 The datasets used in this study can be accessed via zenodo (link to be updated) upon request.

516

517 **FIGURES**

518 **Figure 1: Machine learning workflow for developing and validating an abbreviated smell**  
519 **test.** Details of the workflow are as indicated and described in Methods and Result sections of  
520 the main text. SST-ID = Sniffin' Sticks Identification test. UPSIT = University of Pennsylvania  
521 Smell Identification Test. DeNoPa = De Novo Parkinson Study. PROBE = Prognostic  
522 Biomarkers in Parkinson Disease. HC = healthy control. PD = Parkinson disease. DLB =  
523 dementia with Lewy bodies. MSA = multiple system atrophy. PSP = progressive supranuclear  
524 palsy. ROC = receiver operating characteristic. AUC = area under the ROC curve.

525 **Figure 2: Distribution of olfaction scores using two established tests for different diagnostic**  
526 **groups with parkinsonism in three cohorts.** Cummings estimation plots (a-d) were used to  
527 illustrate and compare smell test score distributions in each diagnostic group: (a) for UPSIT in  
528 the Ottawa Trial cohort, (b) for UPSIT in the PROBE cohort, (c) for SST-ID in the DeNoPa  
529 cohort, (d) UPSIT and SST-ID scores were transformed to percentiles based on age- and sex-  
530 adjusted norms in the combined cohorts. Each data point in the upper panels represents the score  
531 of one participant, and colors represent different groups and diagnosis, as shown in legends. The  
532 vertical lines in the upper panels represent the conventional mean  $\pm$  standard deviation error bars.  
533 The lower panels show the mean group difference (the effect size) and its 95% confidence  
534 interval (CI) estimated by bias-corrected and accelerated bootstrap, using healthy controls as the  
535 reference group. Panels in (e) show ROC curves and AUC values with 95% confidence interval  
536 (CI) for smell tests in each cohort (indicated by different colors; individual scores shown in a-d)  
537 to distinguish PD/DLB versus HC groups (left) and PD/DLB versus MSA/PSP groups (right).  
538 Abbreviations as in **Figure 1**.

539 **Figure 3: Individual scent performances in differentiating PD/DLB from HC groups.** SST-  
540 ID scents are shown using baseline DeNoPa data (a, b) and UPSIT scents for the Ottawa Trial  
541 cohort (c, d). Panels (a) and (c) illustrate distribution of AUC values of each scent across 10-fold  
542 cross-validation using violin plots, with 25%, 50%, and 75% quantile lines. The scents are  
543 ordered in descending order of their mean single-scent AUC value; the color of each scent  
544 changes gradually from the most to the least discriminative value, as indicated by the legend.  
545 Scents shared by both tests are highlighted by bold italic font. Panels (b) and (d) show the

546 percentage of subjects correctly identifying each scent within both groups in each corresponding  
547 cohort. Abbreviations as in **Figure 1**.

548 **Figure 4: Comparison of scent rankings in this study versus previously published ones.**  
549 Panels (a) and (b) show scent rankings of SST-ID and UPSIT, respectively. “This study”  
550 columns show scent rankings from **Figure 3**, and the neighboring columns show corresponding  
551 rankings from other studies, as indicated at the x-axis. The “Average” column of each panel  
552 shows the scent ranking generated by averaging results from 5 separate rankings. Each scent is  
553 represented using the format “index-scent” in the “Average” ranking, and as index only in others.  
554 The lines track how each scent’s rank changes from study to study. Color of each scent changes  
555 gradually from the most to the least discriminative odorant defined by “Average”. Based on  
556 these, 7 best-performing scents in SST-ID (a) and 12 best-performing scents in UPSIT (b) are  
557 tracked by solid lines. Note, rankings by Mahlknecht *et al.* and Morley *et al.* included only the  
558 top 12 scents.

559 **Figure 5: Influence of distractors in multiple-choice smell tests for five shared scents**  
560 **selected.** Panels with odd numbers show the Item Characteristic Curves (ICCs) of five SST-ID  
561 scents: mint, licorice, banana, orange, and lemon. Panels with even numbers show ICCs of the  
562 corresponding UPSIT scents. In each figure, panels on the left show data for PD/DLB patients,  
563 panels on the right for healthy controls. The x-axis reveals transformed score indices (percentage  
564 rank of the respective scores) within the corresponding group. The y-axis shows the probability  
565 of choosing each option at a particular score index. The correct option of each item is highlighted  
566 using thicker, blue curves. Numbers in the color legends represent option indices. The horizontal  
567 dashed lines represent 50% probability. The vertical dashed lines represent five quantiles (5%,  
568 25%, 50%, 75%, and 95%).

569 **Figure 6: Exploration of smaller subsets of scents tested vs. accuracy in group classification**  
570 **of PD/DLB vs HC.** The x-axis shows the number of individual scents used for each subset  
571 examined; colors represent different scent rankings from separate studies, as indicated by the  
572 legends (see also **Figure 4** and **Table 3**). ‘Shared’ denotes scents used in both UPSIT and SST-  
573 ID; Average, all studies combined; This study, rankings derived using baseline DeNoPa and  
574 Ottawa Trail data. Individual points shown in panels (a) and (d) represent internal validation

575 results, averaging across 10 folds. In panels (b), (c) and (e), each point represents the AUC value  
576 of the corresponding subset using (semi-)external validation sets. The black horizontal, dashed  
577 lines indicate AUC values of the corresponding test when viewed in its entirety. Red horizontal,  
578 dashed lines indicate AUC = 0.9 as a predetermined reference line.

579 **Figure 7: Relationships between scent identification performance, diagnosis, sex and age.**  
580 Panel (a) shows the percentage of persons that correctly identified each scent within the healthy  
581 control (HC) group (indicated by light blue region) and the PD/DLB group (indicated by pink  
582 region) in the corresponding cohorts, separated by sex (bar color). Panels in (b) show the  
583 relationship between age (x-axis), diagnostic group (HC in blue; vs PD/DLB in red) and the  
584 percentage of correctly identified scents (y-axis) for each odorant tested (columns) within each  
585 cohort (row), as indicated on the right.



586 TABLES

587 Table 1: Baseline demographic characteristics and smell test scores of adults enrolled in the three cohorts

Variable	DeNoPa					Ottawa Trial					PROBE				
	HC, N = 109 <sup>1</sup>	PD/DLB, N = 129 <sup>1,4</sup>	MSA/PSP, N = 9 <sup>1,4</sup>	p- value <sup>2</sup>	q- value <sup>3</sup>	HC, N = 118 <sup>1</sup>	PD/DLB, N = 70 <sup>1,5</sup>	MSA/PSP, N = 6 <sup>1,5</sup>	p- value <sup>2</sup>	q- value <sup>3</sup>	HC, N = 54 <sup>1</sup>	PD/DLB, N = 102 <sup>1,6</sup>	MSA/PSP, N = 53 <sup>1,6</sup>	p- value <sup>2</sup>	q- value <sup>3</sup>
<b>Sex</b>				0.6	0.7				0.005	0.015				0.058	0.087
Female	42 (39%)	45 (35%)	2 (22%)			74 (63%)	29 (41%)	5 (83%)			28 (52%)	33 (33%)	18 (34%)		
Male	67 (61%)	84 (65%)	7 (78%)			44 (37%)	41 (59%)	1 (17%)			26 (48%)	67 (67%)	35 (66%)		
<b>Age</b>	65 (60, 70)	66 (58, 72)	72 (65, 76)	0.073	0.15	68 (58, 73)	68 (60, 74)	66 (63, 70)	0.9	0.9	59 (55, 69)	55 (55, 69)	67 (61, 75)	0.002	0.007
<b>Parkinsonism duration at baseline in months</b>	NA	14 (9, 24)	12 (6, 33)	0.7	0.7	NA	84 (36, 132)	48 (30, 66)	0.078	0.12	NA	58 (58, 71)	58 (57, 60)	0.2	0.2
<b>Follow-up time in months</b>	120 (120, 120)	120 (72, 120)	72 (48, 120)	0.001	0.004	ND	ND	ND			ND	ND	ND		
<b>Smell test score<sup>7</sup></b>	12 (11, 14)	7 (4, 9)	10 (9, 11)	<0.001	<0.001	32 (29, 35)	17 (13, 22)	29 (26, 31)	<0.001	<0.001	35 (33, 37)	14 (14, 26)	28 (21, 32)	<0.001	<0.001
<b>Smell test percentile<sup>8</sup></b>	50 (25, 75)	4 (4, 10)	25 (18, 50)	<0.001	<0.001	39 (18, 66)	6 (4, 10)	23 (14, 37)	<0.001	<0.001	56 (27, 73)	4 (4, 17)	23 (9, 42)	<0.001	<0.001
<b>Olfaction<sup>9</sup></b>				<0.001	<0.001				<0.001	<0.001				<0.001	<0.001
Normal	94 (86%)	30 (23%)	7 (78%)			93 (79%)	7 (10%)	3 (50%)			50 (93%)	28 (28%)	34 (64%)		
Hyposmia/Anosmia	15 (14%)	99 (77%)	2 (22%)			25 (21%)	63 (90%)	3 (50%)			4 (7%)	74 (72%)	19 (36%)		

<sup>1</sup> n (%); Median (IQR); <sup>2</sup> Fisher's exact test; Kruskal-Wallis rank sum test; <sup>3</sup> False discovery rate correction for multiple testing

<sup>4</sup> PD: n = 126; DLB: n = 3; MSA: n = 4; PSP: n = 5. <sup>5</sup> PD: n = 69; DLB: n = 1; MSA: n = 5; PSP: n = 1. <sup>6</sup> PD: n = 102; MSA: n = 27; PSP: n = 26.

<sup>7</sup> SST-ID scores (0-16) for DeNoPa, UPSIT scores (0-40) for Ottawa Trial and PROBE. <sup>8</sup> Age- and sex-adjusted normalized percentiles.

<sup>9</sup> Hyposmia/anosmia was determined by SST-ID percentile ≤ 10%, and UPSIT percentile ≤ 15%.

588

589 DeNoPa = *De Novo* Parkinson Study. PROBE = Prognostic Biomarkers in Parkinson Disease. IQR = interquartile range. HC = healthy control. PD = Parkinson disease. DLB =

590 dementia with Lewy bodies. MSA = multiple system atrophy. PSP = progressive supranuclear palsy. NA = not applicable. ND = not documented

medRxiv preprint doi: <https://doi.org/10.1101/2024.08.09.24311696>; this version posted August 9, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY 4.0 International license.

591 **Table 2: Discriminative performances of complete smell tests for baseline visits in three**  
592 **cohorts**

Cohort	Test	AUC (95% CI)	Threshold	Sensitivity	Specificity
<b>PD/DLB vs HC</b>					
	SST-ID	0.89 (0.85-0.93)	<= 10.5	0.88	0.81
DeNoPa	SST-TH	0.83 (0.78-0.88)	<= 3.12	0.68	0.9
	SST-DS	0.83 (0.77-0.88)	<= 10.5	0.75	0.78
Ottawa Trial	UPSIT	0.92 (0.88-0.96)	<= 28.5	0.96	0.77
PROBE	UPSIT	0.93 (0.89-0.97)	<= 29.5	0.84	0.89
<b>PD/DLB vs MSA/PSP</b>					
	SST-ID	0.8 (0.69-0.91)	<= 6.5	0.48	1
DeNoPa	SST-TH	0.64 (0.48-0.8)	<= 1.38	0.43	0.89
	SST-DS	0.56 (0.32-0.79)	<= 12.5	0.88	0.33
Ottawa Trial	UPSIT	0.92 (0.85-0.99)	<= 25.5	0.81	1
PROBE	UPSIT	0.69 (0.6-0.78)	<= 26.5	0.76	0.58

593  
594 DeNoPa = *De Novo* Parkinson Study. PROBE = Prognostic Biomarkers in Parkinson Disease.  
595 HC = healthy control. PD = Parkinson disease. DLB = dementia with Lewy bodies. MSA =  
596 multiple system atrophy. PSP = progressive supranuclear palsy. AUC = area under the ROC  
597 curve. CI = confidence interval. SST-ID = Sniffin' Sticks Identification test. SST-TH = Sniffin'  
598 Sticks Threshold test. SST-DS = Sniffin' Sticks Discrimination test. UPSIT = University of  
599 Pennsylvania Smell Identification Test.

600

601 **Table 3: Ranking of scents that are shared by SST-ID and UPSIT kits**

Rank	Scent	Rank in SST-ID <sup>1</sup>	Rank in UPSIT <sup>1</sup>	Average <sup>2</sup>
1	Licorice	2/16	5/40	0.125
2	Banana	4/16	4/40	0.175
3	Clove	8/16	2/40	0.275
4	Rose	7/16	9/40	0.33125
5	Mint	3/16	21/40	0.35625
6	Pineapple	9/16	8/40	0.38125
7	Cinnamon	10/16	16/40	0.5125
8	Lemon	12/16	23/40	0.6625
9	Turpentine	13/16	27/40	0.74375
10	Orange	15/16	29/40	0.83125
11	Leather	11/16	40/40	0.84375

<sup>1</sup> Based on the corresponding "Average" rankings in **Figure 4**.

<sup>2</sup> For each scent, average = (rank in SST-ID + rank in UPSIT) / 2

602

603 SST-ID = Sniffin' Sticks Identification test. UPSIT = University of Pennsylvania Smell

604 Identification Test.

605

606 **Table 4: Performance of the unified abbreviated smell test with 7 scents**

	PD/DLB	HC	MSA/PSP
<b>DeNoPa</b>			
Mean (sd) <sup>1</sup>	3.17 (1.72)	5.69 (1.15)	4.67 (1.00)
AUC (95% CI) <sup>2</sup>		0.88 (0.83-0.92)	0.76 (0.65-0.87)
<b>Ottawa Trial</b>			
Mean (sd) <sup>1</sup>	2.81 (1.56)	5.58 (1.53)	5.00 (1.79)
AUC (95% CI) <sup>2</sup>		0.89 (0.84-0.93)	0.82 (0.62-1)
<b>PROBE</b>			
Mean (sd) <sup>1</sup>	3.49 (1.74)	6.37 (0.90)	4.68 (2.06)
AUC (95% CI) <sup>2</sup>		0.91 (0.87-0.96)	0.68 (0.58-0.77)
<b>Combined</b>			
Mean (sd) <sup>1</sup>	3.20 (1.71)	5.77 (1.31)	4.71 (1.92)
AUC (95% CI) <sup>2</sup>		0.87 (0.85-0.9)	0.73 (0.65-0.8)

<sup>1</sup> For the unified abbreviated smell test with 7 scents, score range is 0-7.

<sup>2</sup> AUC values were for distinguishing patients with PD/DLB from healthy controls or patients with MSA/PSP.

607

608 DeNoPa = *De Novo* Parkinson Study. PROBE = Prognostic Biomarkers in Parkinson Disease.

609 HC = healthy control. PD = Parkinson disease. DLB = dementia with Lewy bodies. MSA =

610 multiple system atrophy. PSP = progressive supranuclear palsy. AUC = area under the ROC

611 curve. CI = confidence interval. sd = standard deviation.

612

613

## 614 **SUPPLEMENTAL METHODS**

### 615 *Detailed description of the study cohorts*

616 *De Novo* Parkinson disease study (DeNoPa): The DeNoPa cohort<sup>[27]</sup> is an ongoing, single-center  
617 study based in Kassel, Germany. The DeNoPa cohort is an observational, longitudinal study of  
618 patients with a newly established diagnosis of PD (UK Brain Bank Criteria<sup>[29]</sup>), who were naïve  
619 to L-DOPA therapy at baseline, and of age- and sex- and education-matched, neurologically  
620 healthy controls (HC). Details of inclusion/exclusion criteria have been described elsewhere.<sup>[27]</sup>  
621 Diagnostic accuracy was ensured by ongoing follow-up visits every two years (as of 2023, 10-  
622 year follow up visits were underway). Data used were received on May 16<sup>th</sup>, 2023.

623 Ottawa (PREDIGT) Trial: The Ottawa Trial is a pilot study to evaluate the performance of a 2-  
624 step screening tool that combines the PREDIGT questionnaire<sup>[8],[9]</sup> and the UPSIT test to  
625 distinguish patients with PD/DLB from age-matched neurologically healthy controls and patients  
626 with various other neurological diseases. Enrolment and assessment of this cross-sectional, case-  
627 control study was completed in March 2024. A manuscript that describes this cohort is in  
628 preparation. Diagnostic accuracy was ensured by independent chart review by three subspecialty-  
629 trained neurologists (JS, MGS, and AF) according to UK Brain Bank Criteria and MDS  
630 Criteria.<sup>[30]</sup>

631 Prognostic Biomarkers in Parkinson Disease (PROBE): PROBE<sup>[28]</sup> is a longitudinal, case-control  
632 study to test biomarkers in PD subjects and controls to determine their feasibility and potential  
633 utility as markers of risk and prognosis for PD. Details of inclusion/exclusion criteria have been  
634 described elsewhere.<sup>[28]</sup> Participants were enrolled from August, 2007 to December, 2008.  
635 Diagnosis of PD, probable MSA, and probable PSP met UK Brain Bank criteria, Consensus  
636 Criteria, and NINDS-PSP Criteria, respectively.

### 637 *Supplemental methods*

638 *10-fold cross validation (CV)*: For each smell test, the discovery dataset was randomly  
639 partitioned into 10 parts, where the case-control ratio was maintained in each part. In each fold,  
640 9/10 parts were used as the development set and the remaining 1 part was used for internal

641 validation. This procedure was repeated for 10 times, and results were showed either in  
642 distribution or average across 10 folds.

643 *Relationship between scent identification and age:* Within each cohort, participants' ages were  
644 segregated into four bins with similar sample size. For each scent, the proportion of correctly  
645 identification was calculated for each bin, and spline smoothing was then used to represent the  
646 relationship between the proportions and age.

647 *Item characteristic curves (ICCs):* The version of ICCs used in this study differed from  
648 traditional parametric ICCs in two aspects: 1) the x-axis was the score percentage rank in [0,100],  
649 not the latent trait on the whole real line; and 2) ICCs represented spline smoothing lines that fit  
650 response data, rather than being fitted to any pre-defined parametric model.<sup>[35]</sup>

651

## 652 SUPPLEMENTAL FIGURES

653 **Supplemental Figure 1: Distribution of Sniffin' Sticks Threshold (SST-TH) and**  
654 **Discrimination (SST-DS) scores for each subject group in the DeNoPa Study at baseline**  
655 **and their ROC curves for group classification.** The Cummings estimation plots (a, b) were  
656 used to illustrate and compare smell test score distributions in each group: (a) SST-TH, (b) SST-  
657 DS. Each data point in the upper panels represents the score of one participant, and colors  
658 represent different groups and diagnosis as shown in the legend. The vertical lines in the upper  
659 panels represent the conventional mean  $\pm$  standard deviation error bars. The lower panels show  
660 the mean group difference (the effect size) and its 95% confidence interval (CI) estimated by  
661 bias-corrected and accelerated bootstrap, using HC as the reference group. Panel (c) shows the  
662 ROC curves of each smell test (indicated by color) to distinguish PD/DLB versus HC (left) and  
663 PD versus OND (right). HC = healthy control. PD = Parkinson disease. DLB = dementia with  
664 Lewy bodies. MSA = multiple system atrophy. PSP = progressive supranuclear palsy. ROC =  
665 receiver operating characteristic. AUC = area under the ROC curve.

666 **Supplemental Figure 2: Percentage differences of correct scent identification between HC**  
667 **and PD/DLB groups (% HC - % PD/DLB) in the DeNoPa (a) and Ottawa Trial (b) cohorts.**  
668 The scents are ordered in descending orders of their mean single-scent AUC value (see Figure 3  
669 (a) and (c)); the color of each scent changes gradually from the most discriminative to the least  
670 discriminative odorant, as indicated by the legend.

671 **Supplemental Figure 3: Influence of distractors in the multiple-choice smell tests: the**  
672 **remaining 6 scents shared by UPSIT and SST-ID.** Panels with odd numbers show the Item  
673 Characteristic Curves (ICCs) of six SST-ID scents, and panels with even numbers show the ICCs  
674 of the corresponding UPSIT scents. In each figure, the left panels are the ICCs using data of the  
675 PD/DLB patients, and the right panels are corresponding to healthy controls. The x-axis is  
676 transformed score indices (percentage rank of the SST-ID score) within the corresponding group.  
677 The y-axis is the probability of choosing each option at a particular score index. The correct  
678 option of each item is highlighted using the thicker blue curves. Numbers in the color legends are  
679 the option indices. The horizontal dashed lines represent 50% probability. The vertical dashed  
680 lines represent five quantiles (5%, 25%, 50%, 75%, and 95%).

681 **Supplemental Figure 4: Item Characteristic Curves (ICCs) of the remaining SST-ID scents**  
682 **using baseline data from the DeNoPa Cohort.** In each panel, the left panels are the ICCs using  
683 data of the PD/DLB patients, and the right panels are corresponding to healthy controls. The x-  
684 axis is transformed score indices (percentage rank of the SST-ID score) within PD and HC group,  
685 respectively. The y-axis is the probability of choosing each option at a particular score index.  
686 The correct option of each item is highlighted using the thicker blue curves. Numbers in the  
687 color legends are the option indices. The horizontal dashed lines represent 50% probability. The  
688 vertical dashed lines represent five quantiles (5%, 25%, 50%, 75%, and 95%).

689 **Supplemental Figure 5: Item Characteristic Curves (ICCs) of the remaining UPSIT scents**  
690 **using the Ottawa Trial cohort.** In each panel, the left panels are the ICCs using data of the  
691 PD/DLB patients, and the right panels are corresponding to healthy controls. The x-axis is  
692 transformed score indices (percentage rank of the SST-ID score) within PD and HC group,  
693 respectively. The y-axis is the probability of choosing each option at a particular score index.  
694 The correct option of each item is highlighted using the thicker blue curves. Numbers in the  
695 color legends are the option indices. The horizontal dashed lines represent 50% probability. The  
696 vertical dashed lines represent five quantiles (5%, 25%, 50%, 75%, and 95%).

697 **Supplemental Figure 6: Range of performances for UPSIT scents in differentiating**  
698 **PD/DLB from MSA/PSP subjects in the PROBE cohort and AUC values for numerical**  
699 **subsets of odorants in group classification.** Panel (a) illustrates distribution of AUC values of  
700 each scent across 10-fold cross-validation using violin plots, with 25%, 50%, and 75% quantile  
701 lines. The scents are ordered in descending order of their mean single-scent AUC value. The  
702 color of each scent changes gradually from the most to the least discriminative odorant, as  
703 indicated by the legend. Panel (b) shows the percentage of subjects correctly identifying each  
704 scent within the MSA/PSP and PD/DLB groups, and panel (c) shows the percentage differences  
705 between the two groups. Scents in panels (b) and (c) follow the same rank order as in panel (a).  
706 In panel (d), the x-axis is the number of scents included for each subset, with individual points  
707 representing average AUC values for the validation set across ten CV folds gathered in PROBE.  
708 The black horizontal, dashed line indicates the corresponding AUC values of the whole test (= 40  
709 scents). The red horizontal, dashed line indicates AUC = 0.8 as a predetermined reference line.



## 710 SUPPLEMENTAL TABLES

## 711 Supplemental Table 1: STARD checklist.

Section & Topic	No	Item	Reported on page #
<b>TITLE OR ABSTRACT</b>			
	1	Identification as a study of diagnostic accuracy using at least one measure of accuracy (such as sensitivity, specificity, predictive values, or AUC)	P1
	2	Structured summary of title, objectives, methods, results and conclusions (for specific guidance, see STARD for Abstracts)	P3
<b>INTRODUCTION</b>			
	3	Scientific and clinical background, including the intended use and clinical role of the index test	Research in context: P4-5 Introduction: P6-7
	4	Study objectives and hypotheses	P7
<b>METHODS</b>			
<i>Study design</i>	5	Whether data collection was planned before the index test and reference standard were performed (prospective study) or after (retrospective study)	P7
<i>Participants</i>	6	Eligibility criteria	P7, Supplemental methods: P2
	7	On what basis potentially eligible participants were identified (such as symptoms, results from previous tests, inclusion in registry)	P7, Supplemental methods: P2
	8	Where and when potentially eligible participants were identified (setting, location and dates)	Supplemental methods: P2
	9	Whether participants formed a consecutive, random or convenience series	P7, Supplemental methods: P2
<i>Test methods</i>	10a	Index test, in sufficient detail to allow replication	Description of SST and UPSIT: P8 The abbreviated smell test: P13
	10b	Reference standard, in sufficient detail to allow replication	Supplemental methods: P2
	11	Rationale for choosing the reference standard (if alternatives exist)	Supplemental methods: P2
	12a	Definition of and rationale for test positivity cut-offs or result categories of the index test, distinguishing pre-specified from exploratory	P9
	12b	Definition of and rationale for test positivity cut-offs or result categories of the reference standard, distinguishing pre-specified from exploratory	Supplemental methods: P2
	13a	Whether clinical information and reference standard results were available to the performers/readers of the index test	P7
	13b	Whether clinical information and index test results were available to the assessors of the reference standard	P7
<i>Analysis</i>	14	Methods for estimating or comparing measures of diagnostic accuracy	P8-10
	15	How indeterminate index test or reference standard results were handled	Supplemental methods: P2
	16	How missing data on the index test and reference standard were handled	Index test: P8. No missing data on the reference standard.
	17	Any analyses of variability in diagnostic accuracy, distinguishing pre-specified from exploratory	P9
	18	Intended sample size and how it was determined	P7
<b>RESULTS</b>			
<i>Participants</i>	19	Flow of participants, using a diagram	P7, Supplemental methods: P2
	20	Baseline demographic and clinical characteristics of participants	Table 1
	21a	Distribution of severity of disease in those with the target condition	P7, Table 1
	21b	Distribution of alternative diagnoses in those without the target condition	Table 1
	22	Time interval and any clinical interventions between index test and reference standard	P7, Supplemental methods: P2
<i>Test results</i>	23	Cross tabulation of the index test results (or their distribution) by the results of the reference standard	Figure 2; Supplemental Figure 1; Tables 1, 2, 4; P10-11, P13-14
	24	Estimates of diagnostic accuracy and their precision (such as 95% confidence intervals)	Figures 2, 6; Supplemental Figures 1, 6; Tables 2, 4; P10-11, P13-14
	25	Any adverse events from performing the index test or the reference standard	The smell tests are non-invasive and there was no adverse event from performing them. P16
<b>DISCUSSION</b>			
	26	Study limitations, including sources of potential bias, statistical uncertainty, and generalisability	P17-18
	27	Implications for practice, including the intended use and clinical role of the index test	P5, P15-17

OTHER INFORMATION		
28	Registration number and name of registry	Supplemental methods: P2
29	Where the full study protocol can be accessed	Supplemental methods: P2
30	Sources of funding and other support; role of funders	P4, P19

medRxiv preprint doi: <https://doi.org/10.1101/2024.08.09.24311696>; this version posted August 9, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](#).

712 Supplemental Table 2: Comparison of cohorts and methods of this study and eight published studies.

	UPSIT					SST-ID				
	<i>This study</i>	<i>Bohnen et al.</i>	<i>Hawkes et al.</i>	<i>Joseph et al.</i> <sup>8</sup>	<i>Morley et al.</i>	<i>This study</i>	<i>Boesveldt et al.</i>	<i>Casjens et al.</i>	<i>Lo et al.</i> <sup>9</sup>	<i>Mahlknecht et al.</i>
<b>Study cohorts</b>										
<b>Discovery</b>										
Cohort name/location	Ottawa Trial	University of Pittsburgh	Ipswich	PREDICT-PD	Michael J. Crescenz VA Medical Center in Philadelphia and the University of Pennsylvania	DeNoPa	VUMC; LUMC	ParkCHIP	Oxford Discovery	Innsbruck and Bruneck
PD	<b>70 (37.2)</b> <sup>1,6</sup>	<b>27 (50)</b> <sup>1</sup>	<b>96 (50)</b> <sup>1</sup>	<b>40 (4.3)</b> <sup>1</sup>	<b>314 (50)</b> <sup>1</sup>	<b>129 (54.2)</b> <sup>1,7</sup>	<b>404 (72.9)</b> <sup>1</sup>	<b>148 (50)</b> <sup>1</sup>	<b>890 (74)</b> <sup>1</sup>	<b>134 (28.5)</b> <sup>1</sup>
Male	41 (59) <sup>1</sup>	20 (74.1) <sup>1</sup>	49 (51) <sup>1</sup>	30 (75) <sup>1</sup>	261 (83) <sup>1</sup>	84 (65) <sup>1</sup>	253 (62.6) <sup>1</sup>	78 (52.7) <sup>1</sup>	569 (64) <sup>1</sup>	84 (62.7) <sup>1</sup>
Age in years	68 (60, 74) <sup>2</sup>	60 (11.1) <sup>3</sup>	57 (27-81) <sup>4</sup>	63.8 (9.6) <sup>3</sup>	67.4 (10.0) <sup>3</sup>	66 (58, 72) <sup>2</sup>	61.5 (40-90) <sup>4</sup>	67 (14) <sup>2</sup>	66.5 (9.6) <sup>3</sup>	68.0 (8.8) <sup>3</sup>
PD duration in years	7 (3, 11) <sup>2</sup>	2.5 (2.7) <sup>3</sup>	unknown	unknown	unknown	1.2 (0.75, 2) <sup>2</sup>	0-44 <sup>5</sup>	n.a.	1.2 (0.9) <sup>3</sup>	6.2 (4.8) <sup>3</sup>
HC	<b>118 (62.8)</b> <sup>1</sup>	<b>27 (50)</b> <sup>1</sup>	<b>96 (50)</b> <sup>1</sup>	<b>891 (95.7)</b> <sup>1</sup>	<b>314 (50)</b> <sup>1</sup>	<b>109 (45.8)</b> <sup>1</sup>	<b>150 (27.1)</b> <sup>1</sup>	<b>148 (50)</b> <sup>1</sup>	<b>313 (26)</b> <sup>1</sup>	<b>336 (71.5)</b> <sup>1</sup>
Male	44 (37) <sup>1</sup>	20 (74.1) <sup>1</sup>	39 (40.6) <sup>1</sup>	343 (38.5) <sup>1</sup>	261 (83) <sup>1</sup>	67 (61) <sup>1</sup>	87 (58) <sup>1</sup>	81 (54.7) <sup>1</sup>	165 (53) <sup>1</sup>	156 (46.4) <sup>1</sup>
Age in years	68 (58, 73) <sup>2</sup>	60 (7) <sup>3</sup>	41.7 (18-78) <sup>4</sup>	67.3 (4.8) <sup>3</sup>	67.4 (10.0) <sup>3</sup>	65 (60, 70) <sup>2</sup>	59.2 (45-78) <sup>4</sup>	62 (16) <sup>2</sup>	64.4 (9.8) <sup>3</sup>	68.8 (8.3) <sup>3</sup>
<b>(Semi-)External Validation</b>	PROBE: PD = 102, HC = 54	n.a.	n.a.	PREDICT-PD: HC = 191, who have completed UPSIT in only year 3.	UCL: PD = 167, HC = 130 Barts: PD = 176, HC = 177	DeNoPa at 48 months: PD=114, HC=101 DeNoPa at 72 months: PD=91, HC=93	n.a.	n.a.	Tracking cohort: 452	VUMC; LUMC: PD = 400, HC = 150 Vienna: PD = 112, HC = 120
<b>Methods</b>										
Internal validation	10-fold cross-validation	n.a.	n.a.	n.a.	n.a.	10-fold cross-validation	n.a.	10-fold cross-validation	Data balance with leave-one-out cross-validation (LOOCV)	n.a.

medRxiv preprint doi: <https://doi.org/10.1101/2024.08.09.24311696>; this version posted August 9, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY 4.0 International license.

Ranking scents	AUC values of each question in distinguishing PD/DLB from HC	Difference of percentages responding correctly in each group (%HC - %PD)	Difference of percentages responding correctly in each group (%HC - %PD)	Difference of percentages responding correctly in each group (%HC - %PD)	40 scents were ranked using 5 methods and only the top-12 of each ranking were reported: 1) the absolute difference in percentage of PD and control subjects answering incorrectly (Difference), 2) odds ratio, 3) discriminant function analysis (Discriminant), 4) logistic regression (Regression), 5) a weighted average combining the first four methods (Combined).	AUC values of each question in distinguishing PD from HC	Difference of percentages responding correctly in each group (%HC - %PD)	Random forest with permutation accuracy	Random forest with predictor importance (Gini diversity index)	L1-regularized logistic regression implementing the least absolute shrinkage and selection operator (the LASSO)
----------------	--	--	--	--	---	--	--	---	--	---

Values are <sup>1</sup>n (%), <sup>2</sup>median (Inter-quartile range (IQR)), <sup>3</sup>mean (standard deviation (SD)), <sup>4</sup>mean (range), <sup>5</sup>range,

<sup>6</sup>Also includes one patient with dementia with Lewy bodies (DLB). <sup>7</sup>Also includes three patients with DLB.

<sup>8</sup>The paper assessed all combinations of 1–7 smells from UPSIT and identified 28 “winning” smell combinations that had highest combined sensitivity and specificity to define hyposmia within 1001 healthy controls in PREDIGT-PD. The ranking below was not directly used to develop the optimal subset of scents.

<sup>9</sup>The discovery cohort was 267 HC with good (normosmia/super-smeller) sense of smell and 721 PD with poor (functional anosmia/hyposmia) sense of smell as defined by age- and sex-specific percentiles: functional anosmia: SST-ID scores ≤ 8; hyposmia: < 10<sup>th</sup> percentile; super-smeller: > 90<sup>th</sup> percentile. The outcome was to classify participants with poor vs normal smell, not to classify PD vs HC, nor conversion to PD in the PD cohort.

medRxiv preprint doi: <https://doi.org/10.1101/2024.08.09.24311696>; this version posted August 9, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY 4.0 International license.

714 **Supplemental Table 3: Relationship between the whole SST-ID score (DeNoPa) and the**  
 715 **whole UPSIT score (Ottawa Trial and PROBE) with age, sex, and diagnostic groups.**

Variable	DeNoPa			Ottawa Trial			PROBE		
	Beta	95% CI <sup>1</sup>	p-value	Beta	95% CI <sup>1</sup>	p-value	Beta	95% CI <sup>1</sup>	p-value
<b>Age</b>	-0.05	-0.09, -0.01	<b>0.022</b>	-0.22	-0.31, -0.14	<b>&lt;0.001</b>	-0.21	-0.31, -0.12	<b>&lt;0.001</b>
<b>Sex</b>									
Female	—	—		—	—		—	—	
Male	-0.66	-1.4, 0.06	0.073	-2.9	-4.6, -1.2	<b>&lt;0.001</b>	-1.9	-3.9, 0.07	0.06
<b>Group</b>									
HC	—	—		—	—		—	—	
PD/DLB	-5	-5.7, -4.3	<b>&lt;0.001</b>	-13	-14, -11	<b>&lt;0.001</b>	-13	-15, -11	<b>&lt;0.001</b>
MSA/PSP	-1.4	-3.3, 0.51	0.2	-1.8	-6.5, 2.9	0.5	-6.4	-9.2, -3.7	<b>&lt;0.001</b>

<sup>1</sup> CI = Confidence Interval

716

717 **REFERENCES**

- 718 [1] Haehner A, Hummel T, Reichmann H. Olfactory loss in Parkinson's disease. *Parkinsons*  
719 *Dis.* 2011;450939.
- 720 [2] Ross GW, Petrovitch H, Abbott RD, et al. Association of olfactory dysfunction with risk  
721 for future Parkinson's disease. *Ann Neurol.* 2008 Feb;63(2):167-73.
- 722 [3] Fereshtehnejad SM, Yao C, Pelletier A, Montplaisir JY, Gagnon JF, Postuma RB.  
723 Evolution of prodromal Parkinson's disease and dementia with Lewy bodies: a prospective  
724 study. *Brain.* 2019 Jul 1;142(7):2051-2067.
- 725 [4] McKinnon JH, Demaerschalk BM, Caviness JN, Wellik KE, Adler CH, Wingerchuk DM.  
726 Sniffing out Parkinson disease: can olfactory testing differentiate parkinsonian disorders?  
727 *Neurologist.* 2007 Nov;13(6):382-5.
- 728 [5] Nalls MA, McLean CY, Rick J, et al. Diagnosis of Parkinson's disease on the basis of  
729 clinical and genetic classification: a population-based modelling study. *Lancet Neurol.* 2015  
730 Oct;14(10):1002-9.
- 731 [6] Bestwick JP, Auger SD, Simonet C, et al. Improving estimation of Parkinson's disease  
732 risk-the enhanced PREDICT-PD algorithm. *NPJ Parkinsons Dis.* 2021 Apr 1;7(1):33.
- 733 [7] Heinzl S, Berg D, Gasser T, Chen H, Yao C, Postuma RB; MDS Task Force on the  
734 Definition of Parkinson's Disease. Update of the MDS research criteria for prodromal  
735 Parkinson's disease. *Mov Disord.* 2019 Oct;34(10):1464-1470.
- 736 [8] Schlossmacher MG, Tomlinson JJ, Santos G, et al. Modelling idiopathic Parkinson  
737 disease as a complex illness can inform incidence rate in healthy adults: the PR EDIGT score.  
738 *Eur J Neurosci.* 2017 Jan;45(1):175-191.
- 739 [9] Li J, Mestre TA, Mollenhauer B, et al. Evaluation of the PREDIGT score's performance  
740 in identifying newly diagnosed Parkinson's patients without motor examination. *NPJ*  
741 *Parkinsons Dis.* 2022 Jul 29;8(1):94.
- 742 [10] Doty RL. Psychophysical testing of smell and taste function, *Handbook of Clinical*  
743 *Neurology*, 164 (2019), pp. 229-246.
- 744 [11] Hummel T, Sekinger B, Wolf SR, Pauli E, Kobal G. 'Sniffin' sticks': olfactory  
745 performance assessed by the combined testing of odor identification, odor discrimination and  
746 olfactory threshold. *Chem Senses.* 1997 Feb;22(1):39-52.

- 747 [12] Rumeau C, Nguyen DT, Jankowski R. How to assess olfactory performance with the  
748 Sniffin' Sticks test(®). *Eur Ann Otorhinolaryngol Head Neck Dis.* 2016 Jun;133(3):203-6.
- 749 [13] Boesveldt S, Verbaan D, Knol DL, et al. (2008) A comparative study of odor  
750 identification and odor discrimination deficits in Parkinson's disease. *Mov Disord* 23: 1984–  
751 1990.
- 752 [14] Hawkes CH, Shephard BC, Daniel SE. Olfactory dysfunction in Parkinson's disease. *J*  
753 *Neurol Neurosurg Psychiatry.* 1997 May;62(5):436-46.
- 754 [15] Bohnen NI, Gedela S, Kuwabara H, et al. Selective hyposmia and nigrostriatal  
755 dopaminergic denervation in Parkinson's disease. *J Neurol.* 2007 Jan;254(1):84-90.
- 756 [16] Morley JF, Cohen A, Silveira-Moriyama L, et al. Optimizing olfactory testing for the  
757 diagnosis of Parkinson's disease: item analysis of the university of Pennsylvania smell  
758 identification test. *NPJ Parkinsons Dis.* 2018 Jan 15;4:2.
- 759 [17] Joseph T, Auger SD, Peress L, et al. Screening performance of abbreviated versions of  
760 the UPSIT smell test. *J Neurol.* 2019 Aug;266(8):1897-1906.
- 761 [18] Casjens S, Eckert A, Woitalla D, et al. Diagnostic value of the impairment of olfaction in  
762 Parkinson's disease. *PLoS One* 2013;8:e64735.
- 763 [19] Mahlknecht P, Pechlaner R, Boesveldt S, et al. Optimizing odor identification testing as  
764 quick and accurate diagnostic tool for Parkinson's disease. *Mov Disord* 2016;31:1408–1413.
- 765 [20] Lo C, Arora S, Ben-Shlomo Y, et al. Olfactory Testing in Parkinson Disease and REM  
766 Behavior Disorder: A Machine Learning Approach. *Neurology.* 2021 Apr 13;96(15):e2016-  
767 e2027.
- 768 [21] Double KL, Rowe DB, Hayes M, et al. Identifying the pattern of olfactory deficits in  
769 Parkinson disease using the brief smell identification test. *Arch Neurol.* 2003 Apr;60(4):545-  
770 9.
- 771 [22] Chou KL, Bohnen NI. Performance on an Alzheimer-selective odor identification test in  
772 patients with Parkinson's disease and its relationship with cerebral dopamine transporter  
773 activity. *Parkinsonism Relat Disord.* 2009 Nov;15(9):640-3.
- 774 [23] Gerkin RC, Adler CH, Hentz JG, et al. Improved diagnosis of Parkinson's disease from a  
775 detailed olfactory phenotype. *Ann Clin Transl Neurol.* 2017 Sep 8;4(10):714-721.



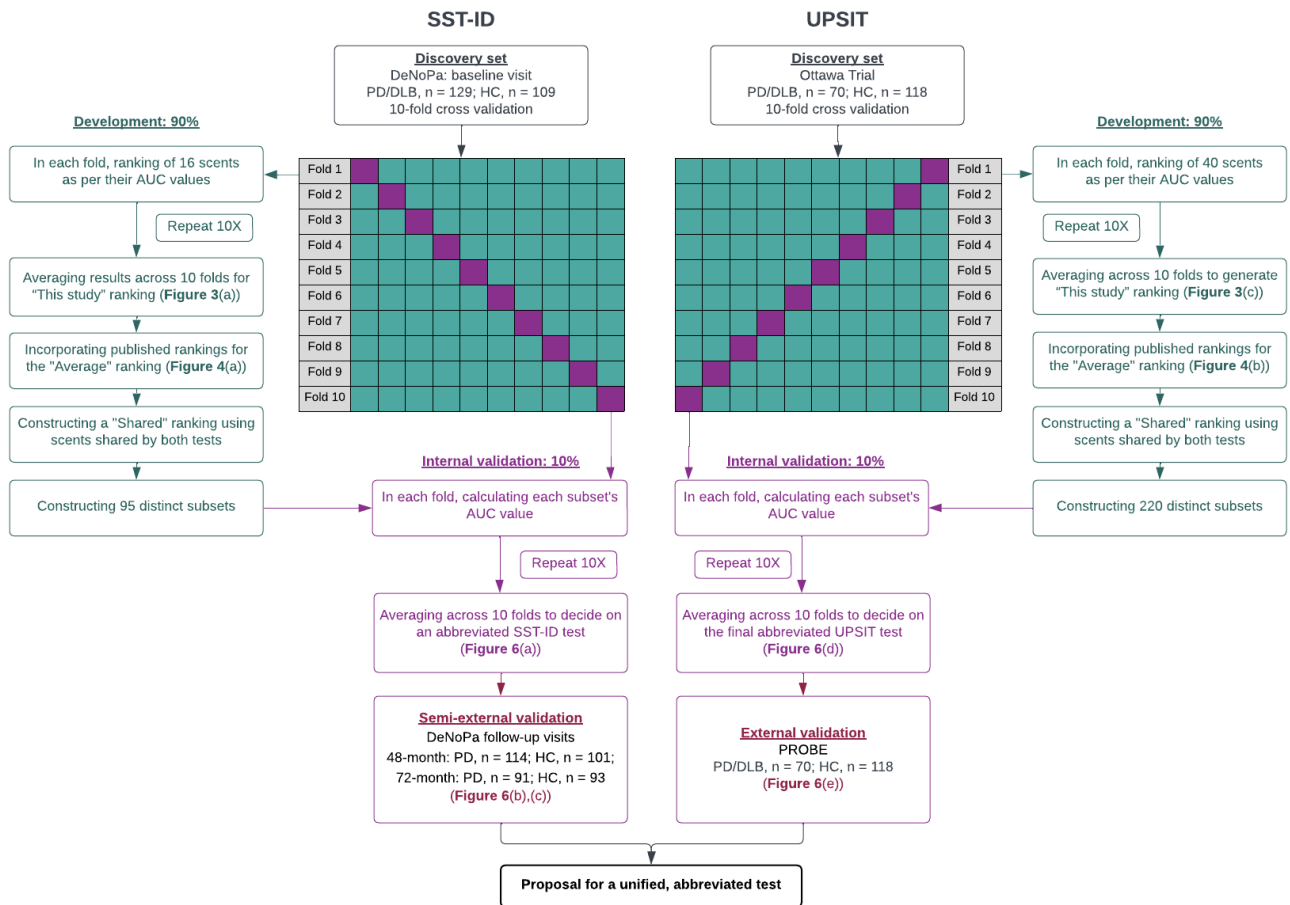
- 776 [24] Auger SD, Kanavou S, Lawton M, et al. Testing Shortened Versions of Smell Tests to  
777 Screen for Hyposmia in Parkinson's Disease. *Mov Disord Clin Pract*. 2020 Mar 21;7(4):394-  
778 398.
- 779 [25] Vaswani PA, Morley JF, Jennings D, Siderowf A, Marek K; PARS Investigators.  
780 Predictive value of abbreviated olfactory tests in prodromal Parkinson disease. *NPJ*  
781 *Parkinsons Dis*. 2023 Jun 29;9(1):103.
- 782 [26] Cohen JF, Korevaar DA, Altman DG, et al. STARD 2015 guidelines for reporting  
783 diagnostic accuracy studies: explanation and elaboration. *BMJ Open*. 2016 Nov 14;6(11):  
784 e012799.
- 785 [27] Mollenhauer B, Trautmann E, Sixel-Döring F, et al. Nonmotor and diagnostic findings in  
786 subjects with de novo Parkinson disease of the DeNoPa cohort. *Neurology*. 2013 Oct  
787 1;81(14):1226-34.
- 788 [28] Diagnostic and Prognostic Biomarkers in Parkinson Disease.  
789 [https://www.ninds.nih.gov/health-information/clinical-trials/diagnostic-and-prognostic-](https://www.ninds.nih.gov/health-information/clinical-trials/diagnostic-and-prognostic-biomarkers-parkinson-disease)  
790 [biomarkers-parkinson-disease](https://www.ninds.nih.gov/health-information/clinical-trials/diagnostic-and-prognostic-biomarkers-parkinson-disease)
- 791 [29] Hughes AJ, Daniel SE, Kilford L, Lees AJ. Accuracy of clinical diagnosis of idiopathic  
792 Parkinson's disease. A clinico-pathological study of 100 cases. *J. Neurol. Neurosurg.*  
793 *Psychiatry* 55, 181–184 (1992).
- 794 [30] Postuma RB, Berg D, Stern M, Poewe W, Olanow CW, Oertel W, et al. MDS clinical  
795 diagnostic criteria for Parkinson's disease. *Mov Disord*. 2015 Oct;30(12):1591-601.
- 796 [31] Cumming G. (2011). Understanding The New Statistics: Effect Sizes, Confidence  
797 Intervals, and Meta-Analysis (1st ed.). Routledge. <https://doi.org/10.4324/9780203807002>
- 798 [32] Hummel T, Kobal G, Gudziol H, Mackay-Sim A. Normative data for the "Sniffin' Sticks"  
799 including tests of odor identification, odor discrimination, and olfactory thresholds: an  
800 upgrade based on a group of more than 3,000 subjects. *Eur Arch Otorhinolaryngol*. 2007  
801 Mar;264(3):237-43.
- 802 [33] Robin X, Turck N, Hainard A, et al (2011). "pROC: an open-source package for R and  
803 S+ to analyze and compare ROC curves." *BMC Bioinformatics*, 12, 77. (version 1.18.5)
- 804 [34] Youden WJ. Index for rating diagnostic tests. *Cancer*. 1950 Jan;3(1):32-5.
- 805 [35] Ramsay JO, Wiberg M, Li J. (2020). Full Information Optimal Scoring. *Journal of*  
806 *Educational and Behavioral Statistics*, 45(3), 297–315.



- 807 [36] Ho JW, Tumkaya T, (2020). *dabestr: Data Analysis using Bootstrap-Coupled Estimation*.  
808 <https://cran.r-project.org/web/packages/dabestr/index.html> (version 0.3.0)
- 809 [37] Ramsay JO, Li J, Wiberg M, Wallmark J, Graves S. *TestGardener: Optimal Analysis of*  
810 *Test and Rating Scale Data*. <https://cran.r-project.org/web/packages/TestGardener/index.html>  
811 (version 3.2.6)
- 812 [38] Wickham H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New  
813 York. ISBN 978-3-319-24277-4, <https://ggplot2.tidyverse.org/>. (version 3.4.4)
- 814 [39] Ramsay JO, Hooker G, Graves S. *fda: Functional Data Analysis*. [https://cran.r-](https://cran.r-project.org/web/packages/fda/index.html)  
815 [project.org/web/packages/fda/index.html](https://cran.r-project.org/web/packages/fda/index.html) (version 6.1.4)
- 816 [40] Lang AE, Siderowf AD, Macklin EA, et al. Trial of Cinpanemab in Early Parkinson's  
817 Disease. *N Engl J Med*. 2022 Aug 4;387(5):408-420.
- 818 [41] Pagano G, Taylor KI, Anzures-Cabrera J, et al. Trial of Prasinezumab in Early-Stage  
819 Parkinson's Disease. *N Engl J Med*. 2022 Aug 4;387(5):421-432.
- 820 [42] Jensen PH, Schlossmacher MG, Stefanis L. Who Ever Said It Would Be Easy? Reflecting  
821 on Two Clinical Trials Targeting  $\alpha$ -Synuclein. *Mov Disord*. 2023 Mar;38(3):378-384.
- 822 [43] Jennings D, Siderowf A, Stern M, et al. (2014) Imaging prodromal Parkinson disease: the  
823 Parkinson Associated Risk Syndrome study. *Neurology* **83**:1739-1746.
- 824 [44] The Lancet. What next in Parkinson's disease? *Lancet*. 2024 Jan 20;403(10423):219.
- 825 [45] Mollenhauer B, Li J, Schlossmacher MG (2023). Persistent Hyposmia as Surrogate for  $\alpha$ -  
826 Synuclein-Linked Brain Pathology. *medRxiv*  
827 <https://www.medrxiv.org/content/10.1101/2023.12.19.23300164v2>
- 828 [46] Stefani A, Iranzo A, Holzkecht E, et al. Alpha-synuclein seeds in olfactory mucosa of  
829 patients with isolated REM sleep behaviour disorder. *Brain*. 2021 May 7;144(4):1118-1126.
- 830 [47] Tomlinson JJ, Shutinoski B, Dong L, et al. Holocranohistochemistry enables the  
831 visualization of  $\alpha$ -synuclein expression in the murine olfactory system and discovery of its  
832 systemic anti-microbial effects. *J Neural Transm (Vienna)*. 2017 Jun;124(6):721-738.
- 833 [48] Martin-Lopez E, Vidyadhara DJ, Liberia T, et al.  $\alpha$ -Synuclein Pathology and Reduced  
834 Neurogenesis in the Olfactory System Affect Olfaction in a Mouse Model of Parkinson's  
835 Disease. *J Neurosci*. 2023 Feb 8;43(6):1051-1071.

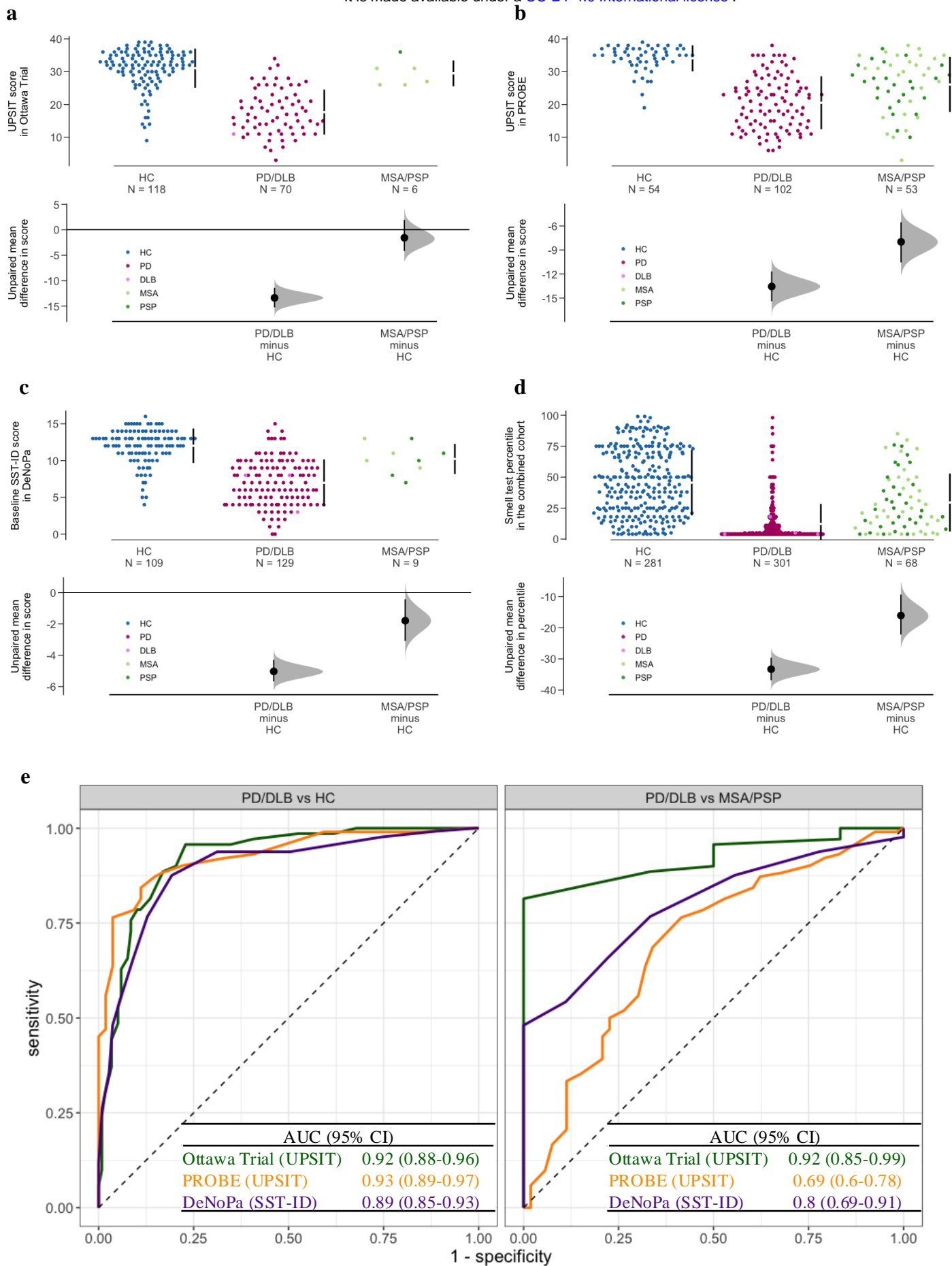
*Li et al., Simplified Olfaction Testing to Identify Patients with Parkinson's*

- 836 [49]Chen F, Liu W, Liu P, *et al.*  $\alpha$ -Synuclein aggregation in the olfactory bulb induces  
837 olfactory deficits by perturbing granule cells and granular-mitral synaptic transmission. *NPJ*  
838 *Parkinsons Dis.* 2021 Dec 13;7(1):114.
- 839 [50]Petit GH, Berkovich E, Hickery M, *et al.* Rasagiline ameliorates olfactory deficits in an  
840 alpha-synuclein mouse model of Parkinson's disease. *PLoS One.* 2013;8(4):e60691.
- 841 [51]Fleming SM, Tetreault NA, Mulligan CK, Hutson CB, Masliah E, Chesselet MF.  
842 Olfactory deficits in mice overexpressing human wildtype alpha-synuclein. *Eur J Neurosci.*  
843 2008 Jul;28(2):247-56.



**Figure 1: Machine learning workflow for developing and validating an abbreviated smell test.**

Details of the workflow are as indicated and described in Methods and Result sections of the main text. SST-ID = Sniffin' Sticks Identification test. UPSIT = University of Pennsylvania Smell Identification Test. DeNoPa = *De Novo* Parkinson Study. PROBE = Prognostic Biomarkers in Parkinson Disease. HC = healthy control. PD = Parkinson disease. DLB = dementia with Lewy bodies. MSA = multiple system atrophy. PSP = progressive supranuclear palsy. ROC = receiver operating characteristic. AUC = area under the ROC curve.



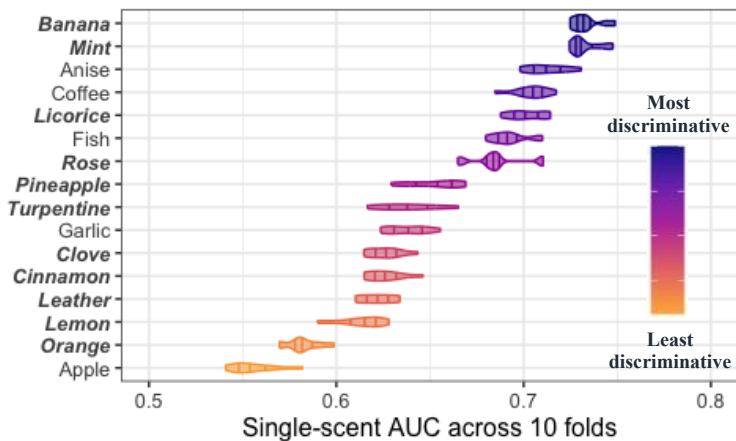
**Figure 2: Distribution of olfaction scores using two established tests for different diagnostic groups with parkinsonism in three cohorts.**

For figure caption, see the next slide.

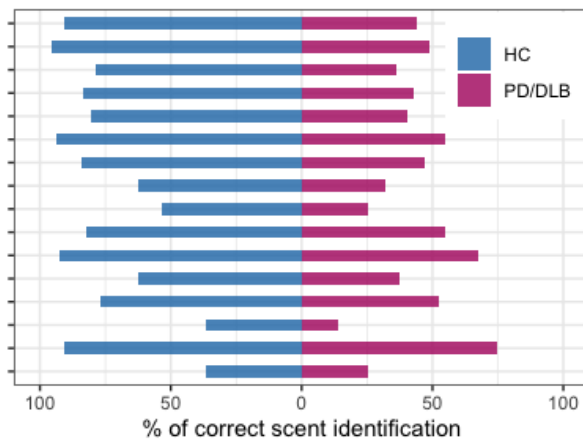
**Figure 2: Distribution of olfaction scores using two established tests for different diagnostic groups with parkinsonism in three cohorts.**

Cummings estimation plots (a-d) were used to illustrate and compare smell test score distributions in each diagnostic group: (a) for UPSIT in the Ottawa Trial cohort, (b) for UPSIT in the PROBE cohort, (c) for SST-ID in the DeNoPa cohort, (d) UPSIT and SST-ID scores were transformed to percentiles based on age- and sex-adjusted norms in the combined cohorts. Each data point in the upper panels represents the score of one participant, and colors represent different groups and diagnosis, as shown in legends. The vertical lines in the upper panels represent the conventional mean  $\pm$  standard deviation error bars. The lower panels show the mean group difference (the effect size) and its 95% confidence interval (CI) estimated by bias-corrected and accelerated bootstrap, using healthy controls as the reference group. Panels in (e) show ROC curves and AUC values with 95% confidence interval (CI) for smell tests in each cohort (indicated by different colors; individual scores shown in a-d) to distinguish PD/DLB versus HC groups (left) and PD/DLB versus MSA/PSP groups (right). Abbreviations as in Figure 1.

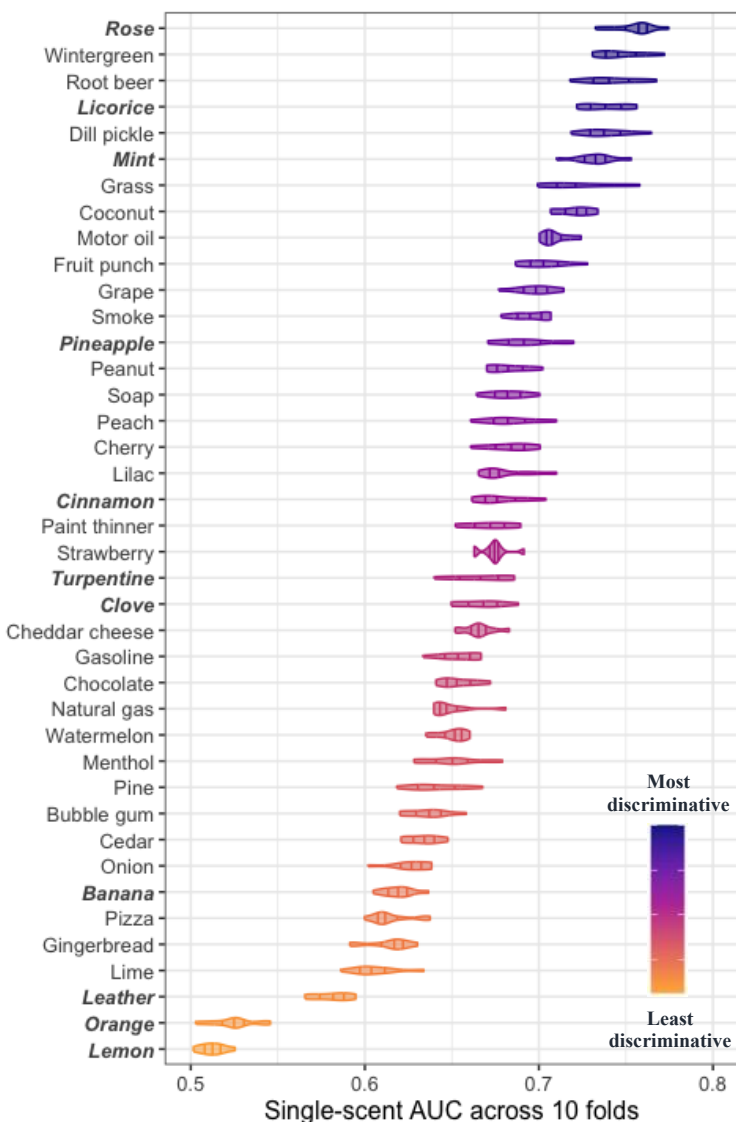
**a**



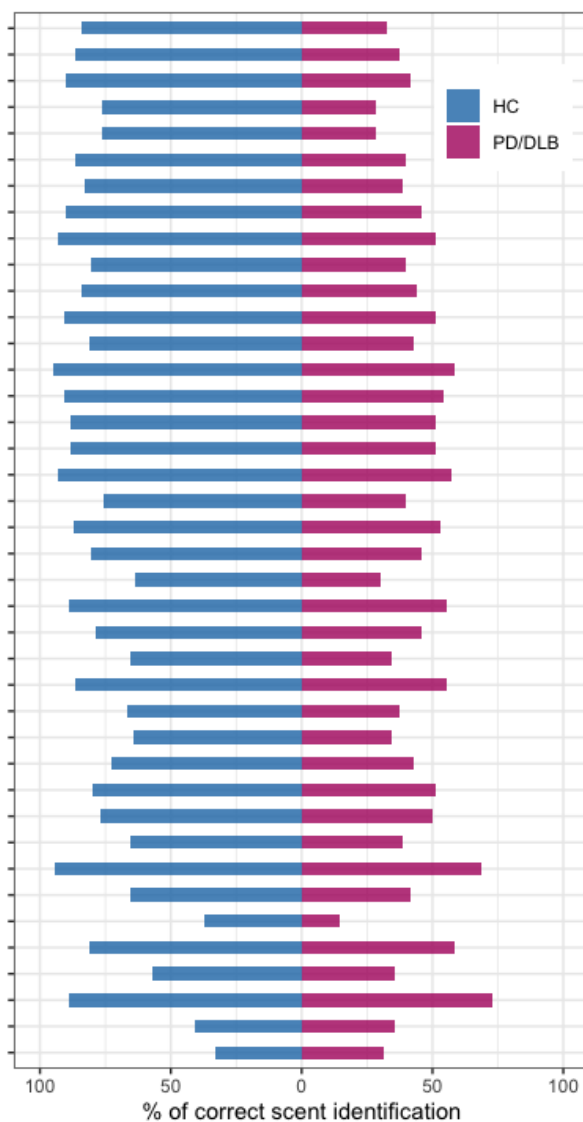
**b**



**c**



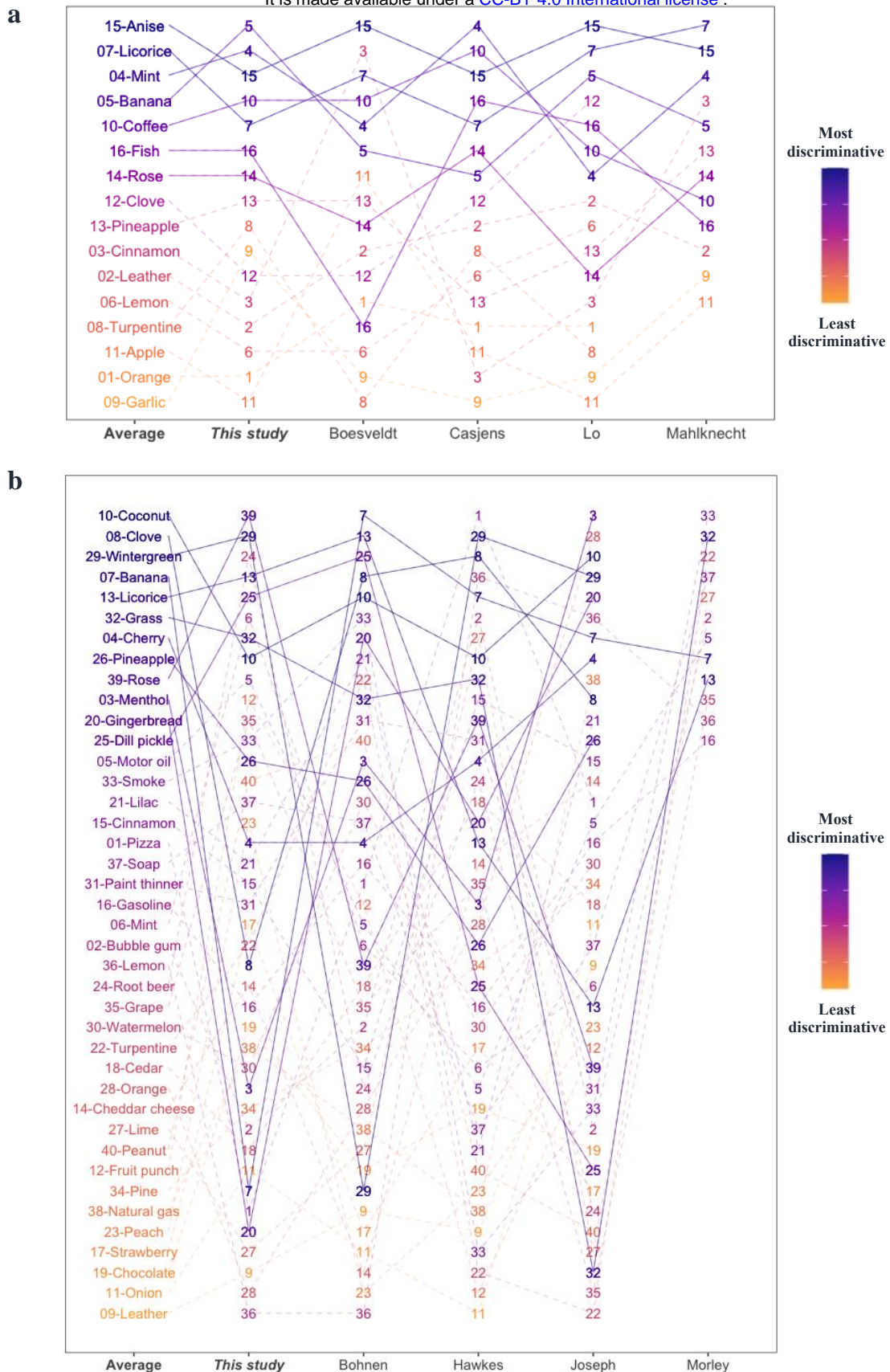
**d**



**Figure 3: Individual scent performances in differentiating PD/DLB from HC groups.**

SST-ID scents are shown using baseline DeNoPa data (a, b) and UPSIT scents for the Ottawa Trial cohort (c, d). Panels (a) and (c) illustrate distribution of AUC values of each scent across 10-fold cross-validation using violin plots, with 25%, 50%, and 75% quantile lines. The scents are ordered in descending order of their mean single-scent AUC value; the color of each scent changes gradually from the most to the least discriminative value, as indicated by the legend. Scents shared by both tests are highlighted by bold italic font. Panels (b) and (d) show the percentage of subjects correctly identifying each scent within both groups in each corresponding cohort. Abbreviations as in Figure 1.



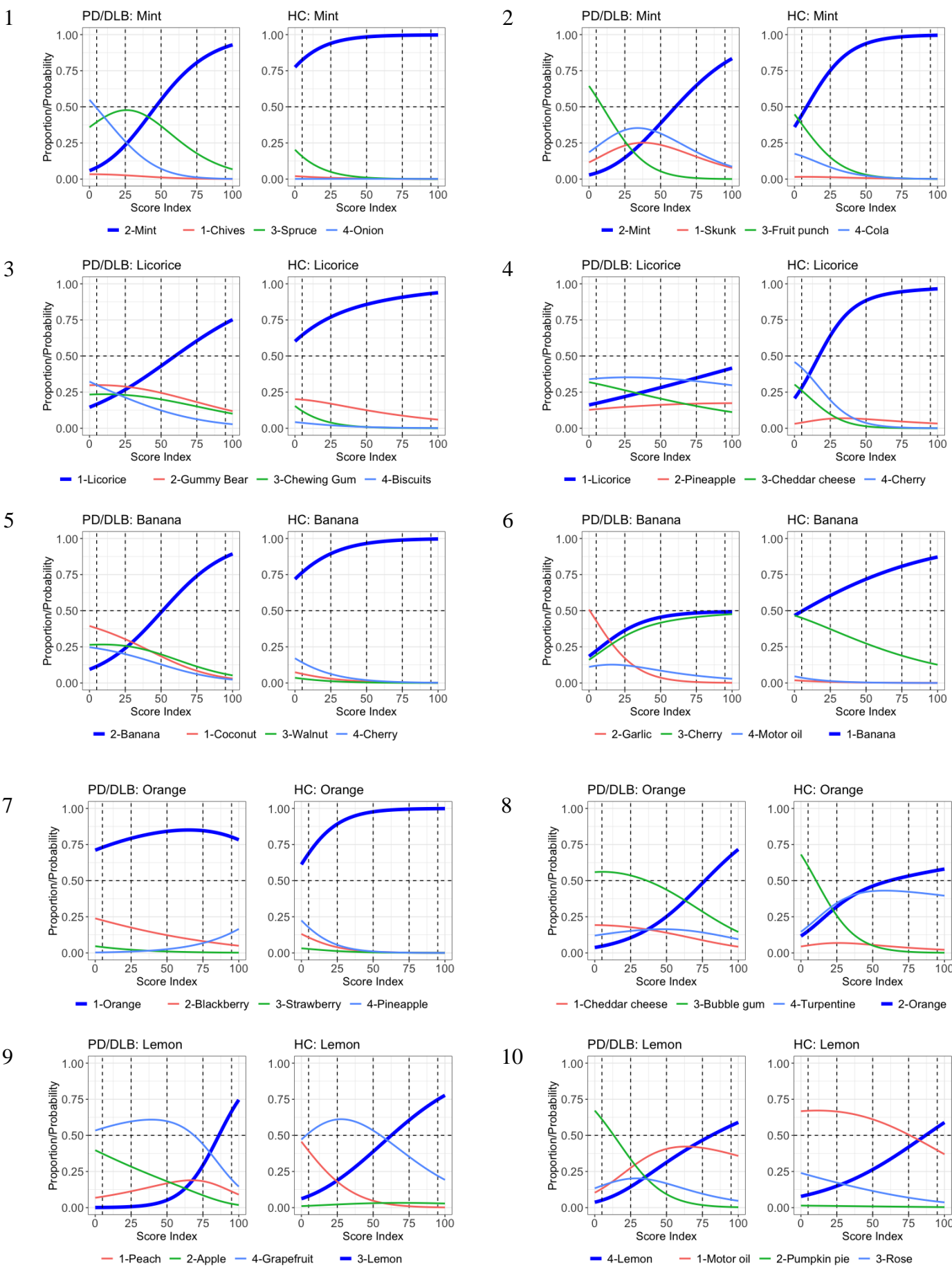


**Figure 4: Comparison of scent rankings in this study versus previously published ones.**

Panels (a) and (b) show scent rankings of SST-ID and UPSIT, respectively. “This study” columns show scent rankings from Figure 3, and the neighboring columns show corresponding rankings from other studies, as indicated at the x-axis. The “Average” column of each panel shows the scent ranking generated by averaging results from 5 separate rankings. Each scent is represented using the format “index-scent” in the “Average” ranking, and as index only in others. The lines track how each scent’s rank changes from study to study. Color of each scent changes gradually from the most to the least discriminative odorant defined by “Average”. Based on these, 7 best-performing scents in SST-ID (a) and 12 best-performing scents in UPSIT (b) are tracked by solid lines. Note, rankings by Mahlkecht *et al.* and Morley *et al.* included only the top 12 scents.

### SST-ID, DeNoPa, baseline

### UPSIT, Ottawa Trial

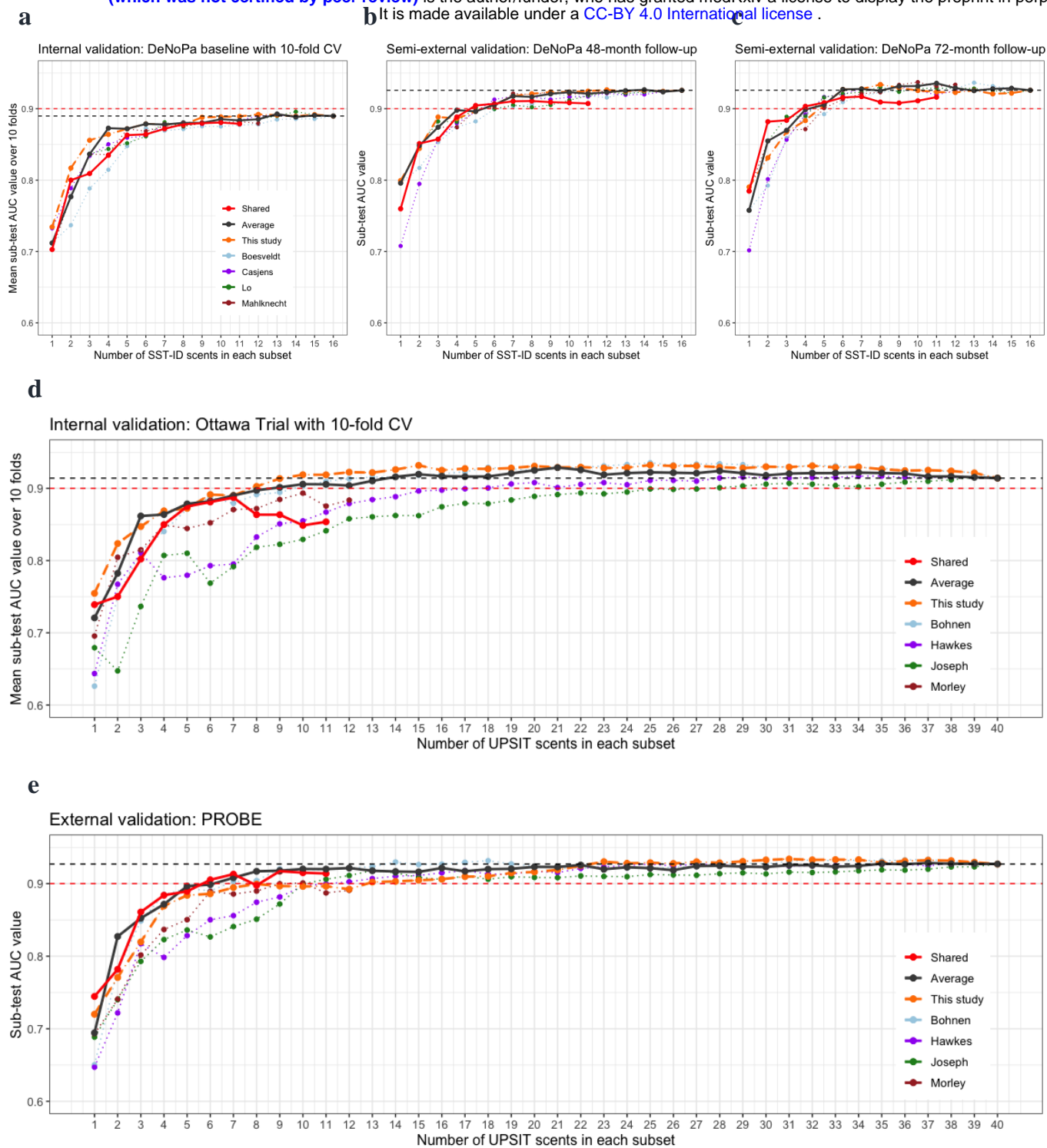


**Figure 5: Influence of distractors in multiple-choice smell tests for five shared scents selected.**

For figure caption, see the next slide.



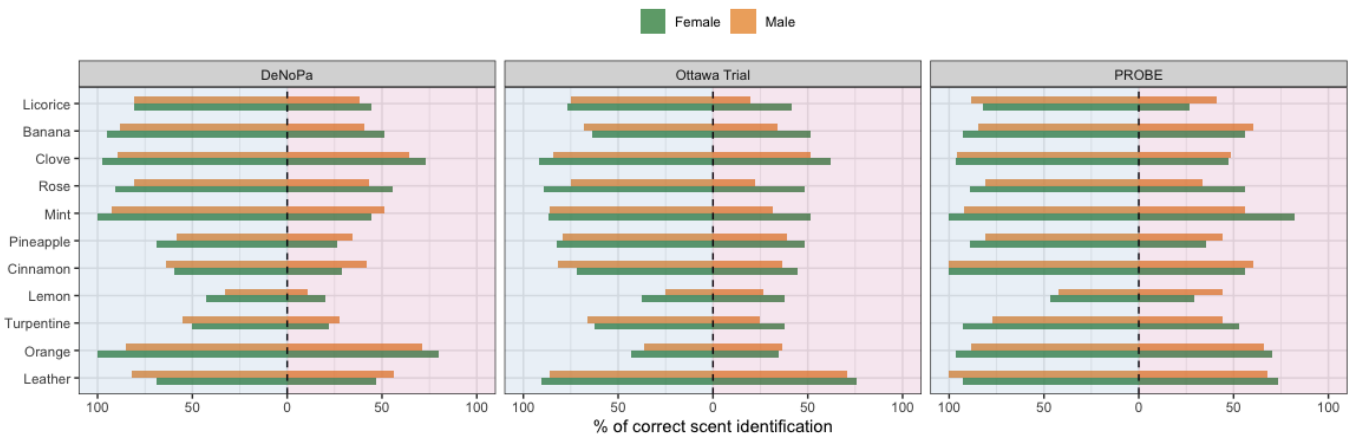
**Figure 5: Influence of distractors in multiple-choice smell tests for five shared scents selected.** Panels with odd numbers show the Item Characteristic Curves (ICCs) of five SST-ID scents: mint, licorice, banana, orange, and lemon. Panels with even numbers show ICCs of the corresponding UPSIT scents. In each figure, panels on the left show data for PD/DLB patients, panels on the right for healthy controls. The x-axis reveals transformed score indices (percentage rank of the respective scores) within the corresponding group. The y-axis shows the probability of choosing each option at a particular score index. The correct option of each item is highlighted using thicker, blue curves. Numbers in the color legends represent option indices. The horizontal dashed lines represent 50% probability. The vertical dashed lines represent five quantiles (5%, 25%, 50%, 75%, and 95%).



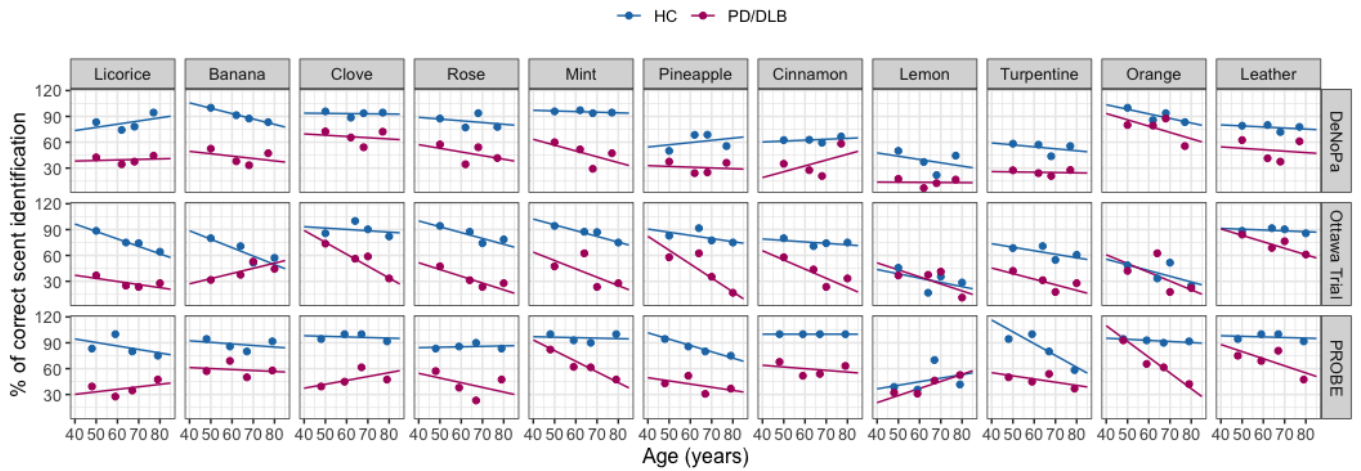
**Figure 6: Exploration of smaller subsets of scents tested vs. accuracy in group classification of PD/DLB vs HC.**

The x-axis shows the number of individual scents used for each subset examined; colors represent different scent rankings from separate studies, as indicated by the legends (see also Figure 4 and Table 3). ‘Shared’ denotes scents used in both UPSIT and SST-ID; Average, all studies combined; This study, rankings derived using baseline DeNoPa and Ottawa Trail data. Individual points shown in panels (a) and (d) represent internal validation results, averaging across 10 folds. In panels (b), (c) and (e), each point represents the AUC value of the corresponding subset using (semi-)external validation sets. The black horizontal, dashed lines indicate AUC values of the corresponding test when viewed in its entirety. Red horizontal, dashed lines indicate AUC = 0.9 as a predetermined reference line.

**a**



**b**



**Figure 7: Relationships between scent identification performance, diagnosis, sex and age.**

Panel (a) shows the percentage of persons that correctly identified each scent within the healthy control (HC) group (indicated by light blue region) and the PD/DLB group (indicated by pink region) in the corresponding cohorts, separated by sex (bar color). Panels in (b) show the relationship between age (x-axis), diagnostic group (HC in blue; vs PD/DLB in red) and the percentage of correctly identified scents (y-axis) for each odorant tested (columns) within each cohort (row), as indicated on the right.