# Smell and Tell: The emergence of olfactory expertise in perfumery students

- 3 <u>Authors:</u> Saive Anne-Lise<sup>1,2</sup>, Plailly Jane<sup>3</sup>, Chambaron Stéphanie<sup>4</sup>, Hourri Youness<sup>1</sup>, Monin
- 4 Carla<sup>1</sup>, Nhouchi Zeineb<sup>5</sup>, Belay Justine<sup>5</sup>, Chalut Pauline<sup>5</sup>, Hanaei Farnaz<sup>5</sup>, Vallet Nadine<sup>5</sup>
- 5 Affiliations
- 6 <sup>1</sup>Lyfe Institute Research and Innovation Center, Écully, France
- 7 <sup>2</sup>Unique Center (Québec Neuro-AI Research Center), Montréal, Québec, Canada
- 8 <sup>3</sup>Centre de Recherche en Neurosciences de Lyon (Inserm U1028 CNRS UMR5292 UCBL)
- 9 CMO Team, Bron, France
- 10 <sup>4</sup> Centre des Sciences du Goût et de l'Alimentation, CNRS, INRAE, Institut Agro, Université
- 11 de Bourgogne, F-21000 Dijon, France.
- <sup>5</sup>Institut Supérieur International du Parfum, de la Cosmétique et de l'Aromatique alimentaire
- 13 (ISIPCA), F-78000 Versailles, France
- 14

# 15 ORCID-IDs:

- 16 Anne-Lise Saive, <u>https://orcid.org/0000-0001-7352-0312</u>
- 17 Jane Plailly, <u>https://orcid.org/0000-0003-4189-1631</u>
- 18 Stéphanie Chambaron, https://orcid.org/
- 19 Zeineb Nhouchi, <u>https://orcid.org/0000-0003-3509-373X</u>
- 20 Nadine Vallet, https://orcid.org/0000-0001-8943-8620
- 21
- 22 Correspondence should be sent to Anne-Lise Saive, <u>alsaive@institutlyfe.com</u>
- 23

# 24 Abstract

- 25 Developing olfactory expertise is essential in professions like perfumery, where the ability to
- 26 describe, categorize, and conceptualize odors is critical. This study investigates how academic
- 27 training during a 1.5-year program at a perfumery school (ISIPCA) shapes olfactory expertise
- 28 of perfumery students. Forty students were assessed at three time points, focusing on odor
- 29 description, evocation, recognition, discrimination, and categorization tasks. Results show
- 30 that training significantly enhanced language abilities related to odor description and
- 31 categorization. Students developed a richer and more precise vocabulary to characterize

32 odors, aligning more closely with expert's terminology and contributing to the formation of a 33 shared olfactory lexicon. Semantic similarity within and between students, as well as with 34 expert references, increased, emphasizing the importance of consistent language use in 35 expertise development. Advanced natural language processing and machine learning tools 36 revealed that the richness of verbal descriptions and semantic similarity were strong 37 predictors of expertise acquisition. In contrast, improvements in non-verbal tasks, such as 38 odor discrimination and recognition, were more limited, suggesting that perceptual abilities 39 may require more extensive training or specialized methods. Building on these results, we 40 propose potential enhancements to olfactory training including reinforced language practice, 41 mental imagery exercises, and sensory discrimination tasks, along with personalized training 42 strategies. These findings highlight the central role of language in the emergence of olfactory 43 expertise and the importance of computational methods for optimizing training programs and 44 advancing educational practices in olfactory science.

Keywords: olfactory expertise, learning, perfume training, natural language processing,
 machine-learning

47

48 **1. Introduction** 

Olfactory expertise is a distinctive and highly specialized skill, cultivated through intensive, deliberate practice and targeted training. This expertise is foundational in professions such as perfumers and sommeliers, where individuals demonstrate extraordinary abilities to recognize, discriminate, and describe a wide array of complex odor profiles. For perfumers, this expertise involves not only mastering the intricate combination of thousands of aromatic compounds to design sophisticated fragrances but also the exceptional ability to mentally construct and evaluate scent compositions without physical stimuli.

56 Achieving olfactory expertise requires both perceptual learning and the development of higher-57 level cognitive processes (for reviews, Parr, 2019; Royet et al., 2013). While olfactory training 58 has been shown to improve sensitivity, these effects are often modest and odor-specific 59 (Haehner et al., 2013; Al Aïn et al., 2019; Zambom-Ferraresi et al., 2021). Olfactory expertise is primarily characterized by enhancements in cognitive functions, including conceptualizing, 60 61 categorizing, and memorizing olfactory information (Solomon, 1997; Hughson and Boakes, 62 2001; Plailly et al., 2012). Short-term training studies also support the role of language in 63 olfactory expertise. Training laypeople to consistently label odors improves their ability to learn 64 and categorize them, resulting in higher accuracy and faster learning (Fournel et al., 2017; 65 Vanek et al., 2021). This suggests a causal link between verbal description and the formation 66 of odor categories, which is particularly noteworthy in the olfactory domain due to the well-67 documented challenges of articulating smells through language (Herz, 2005; Olofsson and 68 Gottfried, 2015). However, while short-term training enhances specific skills like odor labeling, 69 its benefits may not extend to overall olfactory memory or more generalized tasks, as shown by 70 Lesschaeve & Issanchou (1996) and Al Aïn et al. (2019). Together, these findings highlight the 71 nuanced and task-specific nature of olfactory expertise development.

72 In professional perfumers, extensive training leads to superior olfactory mental imagery and 73 richer, more precise verbal descriptions of odors. Plailly et al. (2012) demonstrated that 74 perfumers outperform novices in generating odor mental imagery. This enhanced imagery 75 ability is highly specialized and does not transfer across sensory modalities, suggesting the 76 domain-specific nature of olfactory expertise (Bensafi et al., 2017). Perfumers use more 77 detailed and precise language when describing odors, focusing on chemical and olfactory 78 qualities rather than general hedonic terms, reflecting a shift in semantic odor processing with 79 expertise (Sezille et al., 2014).

80 Similarly, wine experts build their olfactory expertise through improved discrimination, identification, and recognition of wine-related odors. Studies have shown that wine experts have 81 82 a modest yet significant advantage in naming smells and flavors within their domain, 83 attributable to both perceptual and linguistic training (Croijmans and Majid, 2016). Olfactory 84 training has been shown to enhance odor identification abilities, even in novices. For instance, 85 short-term training over five days allowed participants to perform as well as sommeliers in identifying single odorants (Poupon et al., 2018). However, this improvement did not extend to 86 87 mixtures of two or more odorants, highlighting a limitation of short-term training for complex 88 olfactory tasks. In contrast, other research has demonstrated that targeted olfactory training can 89 improve novices' ability to detect individual components within a mixture, suggesting that with 90 focused and potentially longer training, individuals can develop the skills to analyze odor 91 mixtures (Morquecho-Campos et al., 2019). Wine experts also exhibit superior memory for wine odors, a skill that is highly specific to their domain of expertise and reflects the effects of 92 93 specialized training (Parr et al., 2002; Croijmans et al., 2021). Training enhances their mental 94 imagery for wine; the vividness of wine imagery increases only after specialized training, 95 underscoring the role of targeted practice in building olfactory expertise (Croijmans et al., 96 2020). Additionally, wine experts use richer and more concrete language with distinct patterns

when describing wines, which facilitates the translation of nuanced sensory inputs into detailed
verbal descriptions (Sezille et al., 2014; Croijmans et al., 2024).

99 Specialized training in perfumery and wine tasting induces both functional and structural brain 100 changes associated with perceptual and cognitive enhancement. Functional imaging studies 101 reveal that professional perfumers exhibit decreased activation in olfactory and memory regions 102 during odor imagery tasks, indicating enhanced efficiency and reduced cognitive effort (Plailly 103 et al., 2012). Structural brain changes include increased gray matter volume in areas 104 surrounding the olfactory sulcus, such as the gyrus rectus and medial orbital gyrus, with this 105 reorganization positively correlated with experience, even counteracting the effect of age, 106 reflecting the brain's adaptability to olfactory expertise (Delon-Martin et al., 2013). Neural 107 differences are observed in wine experts compared to novices. Wine sommeliers show 108 heightened activation in brain regions responsible for gustatory-olfactory integration and high-109 level cognitive processing (Castriota-Scanderbeg et al., 2005; Pazart et al., 2014; Banks et al., 110 2016; Sreenivasan et al., 2017). Structural adaptations are also evident, with sommeliers 111 exhibiting increased volume in sensory olfactory regions, including the olfactory bulb, right 112 insula, and bilateral entorhinal cortex (Banks et al., 2016; Filiz et al., 2022). Enhanced 113 connectivity between sensory and semantic networks further supports the translation of nuanced 114 sensory inputs into detailed verbal descriptions (Carreiras et al., 2024).

Overall, these findings raise important questions about the extent to which olfactory training generalizes—whether its benefits are limited to specific odor labeling or extend to broader linguistic and cognitive processes associated with olfaction. Furthermore, little is known about the gradual emergence of expertise over time. Although Filiz et al. (2022) investigated how sommelier students develop expertise through a rigorous 1.5-year program involving extensive sensory training in educational settings, most studies still focus primarily on comparisons between experts and novices rather than the process of expertise formation.

122 The present study addresses these gaps by investigating the development of olfactory expertise 123 in perfumery students over a 1.5-year training program. We assessed a wide range of cognitive 124 competencies through tasks including odor categorization, odor description, odor evocation, 125 odor recognition, and odor discrimination. In addition, we performed a detailed analysis of the 126 language used by participants to describe odors using natural language processing (NLP) and 127 machine learning (ML) tools. This research aims to provide a comprehensive understanding of 128 how academic training enhances-or fails to enhance-sensory perception, memory, and 129 related cognitive abilities, contributing to advancements in olfactory education and professional 130 expertise.

131 **2.** Materials and methods

#### 132 **2.1.** Participants

133 The olfactory experiments were conducted at ISIPCA (International School of Perfume, 134 Cosmetics Products, and Food Flavor Formulation) between 2018 and 2021 by ISIPCA faculty 135 and staff, integrating the experiments into the educational curiculum. Since the study was 136 embedded within the school's pedagogical framework and conducted by instructors rather than 137 external researchers, no ethical committee approval was required. The study involved 40 138 ISIPCA students, 35 women and 5 men, with an average age of 17.475 years (SD = 1.585). The 139 gender distribution reflected the demographics of the school's student cohorts, a factor beyond 140 our control.

141 **2.2.Olfactory training** 

Participants were enrolled in the "Specialized Laboratory Technician" diploma, a two-year fulltime program. This curriculum covered foundational knowledge in perfumery, cosmetics, and

food flavors, with olfaction as a major component. The training emphasized olfactory
perception, description, and categorization through structured courses and hands-on sessions.

146 Perfumery Courses - Students studied 120 raw materials over the two-year program. Each raw 147 material was smelled blindly, described verbally, and classified into an olfactory family based 148 on Jean-Noël Jaubert's field of odors (1995). Students assigned three descriptive terms and 149 created detailed sheets for each raw material. For homework, students assessed the volatility of 150 raw materials using scent strips and containers prepared during class. Sessions lasted 3 hours 151 and 30 minutes, including a minimum 20-minute break. The raw materials were organized by 152 olfactory families to facilitate comparative learning, gradually introducing more intense odors 153 as the course progressed. A final exam evaluated students' knowledge. They were required to 154 identify 10 raw materials presented blindly, specifying the olfactory family, descriptors, 155 volatility, and handling characteristics.

156 Food Flavors Courses - The training approach for food flavors followed the same principles 157 as perfumery but included an additional retro-nasal testing step. Raw materials were diluted in 158 water, tasted, and analyzed both ortho and retro-nasally. Students discussed descriptors until 159 reaching a consensus for qualitative and quantitative evaluations. In the first year, sessions 160 covered 6-8 raw materials, increasing to 10 raw materials in the second year. Each session 161 focused on two distinct families (e.g., fruity and toasted). Evaluations progressed from less to 162 more persistent materials. For example, fruity materials with minimal afterglow were evaluated 163 first, followed by more persistent toasted materials. Sessions included one or two breaks. 164 Exams, held two to three times per year, required students to identify 10 raw materials presented 165 blindly, specifying descriptors, families, and persistence characteristics.

#### 166 2.3.Olfactory Stimuli

167 The odorants consisted of monomolecular compounds, essential oils, homemade compositions, 168 and perfumes, all sourced from the olfactory base-sample used in ISIPCA's courses. The 169 monomolecular compounds and essential oils were diluted in alcohol to concentrations 170 appropriate for olfaction, while the perfumes were pure. All stimuli were presented in brown 171 15ml bottles coded with a unique three-digit number. A total of 92 independent samples were 172 selected, each assigned to specific tests: (i) 12 for the evocation and free description tasks, (ii) 173 6 odorants for the triangle test, (iii) 60 for the olfactory recognition test, and (iv) 14 for the 174 categorization test (Table S1).

#### 175 2.4. Visual stimuli in the visual recognition task

176 In the visual recognition task, 120 images were divided into six sets of 20 images each. These 177 images were sourced from previously published datasets, including the DMS48 test (Barbeau 178 et al., 2004). Each participant viewed a unique subset of images across three time points, 179 ensuring that different images were presented in each session. This design minimized 180 familiarity effects and provided a more accurate assessment of recognition memory over time. 181 While unrelated to participants' olfactory training, this task allowed for the assessment of 182 general recognition abilities in a domain not influenced by students' specialized curriculum.

183 **2.5.Experimental tasks** 

184 All tests were carried out in the ISIPCA sensory laboratory. Testing and data collection took 185 place in standard sensory booths, with white lighting and controlled temperature  $(20 \pm 2^{\circ}C)$  and 186 airflow conditions. The expert students participated at three different measurement times: T0 187 (on arrival at the school, before olfactory training started), T1 (6 months after T0), and T2 (18 188 months after T0, 12 months after T1).

189 To assess the effect of academic olfactory training on olfactory capacities, different tasks were190 conducted over four half-days on four separate days at each measurement time:

- **Day 1:** Odor Description, Evocation, and Discrimination tasks
- **Day 2:** Visual and Olfactory Encoding part of the recognition tasks
- **Day 3:** Visual and Olfactory Retrieval part of the recognition tasks
- **Day 4:** Odor Categorization and Categories Description task

195 These tasks were designed to evaluate various aspects of olfactory perception and cognitive 196 processing across the three measurement times and were conducted as described below:

197 **Odor description and evocation:** In the description task, the objective was to generate terms in 198 writing describing the odor by answering the question, "How would you objectively describe 199 this odor?". In the evocation task, students were asked, "What spontaneously comes to mind 200 when you smell this odor?" and were invited to complete the sentences: "This scent makes me 201 want to..." and "This scent reminds me of...". For both tasks, each student was presented with 202 two different olfactory stimuli at each measurement time, with up to two minutes allocated to 203 complete both tasks for each odor. The selection of stimuli was randomized for each student 204 but ensured that student evaluated each odor only once.

Odor discrimination: Students were presented with three strips in a randomized order: two strips contained the same odor, and one strip contained a different odor. They were asked to identify the odd odor. This task involved three trials with the following sets of stimuli: (i) Flowerbomb (\*2) and La Vie Est Belle (\*1), (ii) La Nuit Trésor (\*2) and Black Opium (\*1), and (iii) L'Homme (\*2) and Pink Grapefruit (\*1). The same stimuli were used across all measurement times.

211 Odor recognition: The test was conducted over two consecutive days. On the first day, students
212 were presented with ten target olfactory stimuli, each labeled with a unique three-digit code

213 specific to Day 1 (Encoding) and different from Day 2 (Recognition). They were instructed to 214 smell each stimulus and try to memorize them. On the second day, students were presented with 215 a set of 20 olfactory stimuli, including the 10 target stimuli from the first day and 10 new stimuli 216 (distractors). They were given 5 seconds to smell each stimulus and then indicate whether they 217 had encountered the odor on the first day. In total, 60 odorants were used across all students 218 and time measurements, with different subsets of odorants assigned for each time measurement 219 per student. Responses were recorded using Fizz software during the recognition phase on Day 220 2, with no maximum response time set.

221 Visual recognition: This task followed the same paradigm as the odor recognition task, but with visual stimuli replacing olfactory stimuli. On Day 1 (Encoding), participants were 222 223 presented with a series of pictures using a timed PowerPoint presentation, with each image 224 displayed for 5 seconds. On Day 2 (Recognition), participants' responses were recorded using 225 Fizz software, with no maximum response time set. The inclusion of a visual recognition task 226 served as a control to assess participants' baseline cognitive abilities in a domain unrelated to their training. This distinction allows us to isolate the impact of olfactory training on 227 228 performance and ensure that any differences observed are specific to the sensory modality of 229 interest, rather than reflecting general improvements in recognition or memory processes.

Odor categorization task and categories description: This task aimed at investigating odor learning through categorization. Students were presented with a total of 15 odorants in a twostep test. In the first step, they were asked to group the odors they perceived as harmonious together, forming multiple groups ranging from a minimum of 2 to a maximum of 14. In the subsequent step, students explained their assessments by providing descriptive terms that capture the common attributes of odors within each group. One participant did not complete the task at T2 and therefore was removed from the analyses of this task.

*Odor description and evocation*: The verbal responses from both tasks were first independently 238 239 preprocessed by removing non-alphanumeric characters and converting the text to lowercase 240 before tokenization into individual words and eliminating stop words to normalize the data. The 241 cleaned tokens were then assigned part-of-speech tags (e.g., nouns, adjectives, verbs) using a 242 SpaCy's French language model (fr core news sm). Nouns were analyzed as reprensentations of olfactory sources, while adjectives were identified as descriptors of the olfactory sources. 243 244 The total number of words (including nouns, adjectives, adverbs, and verbs) was also counted 245 as a general proxy for verbal richness. To validate the accuracy of the automated part-of-speech 246 tagging, a manual, exhaustive check was performed, ensuring the correct classification of 247 tokens and refining the data where necessary. After tagging, tokens were lemmatized using 248 SpaCy to reduce words to their root forms, ensuring consistency in lexical analysis. Lexical 249 diversity was then measured using two complementary measures: type-token ratio (TTR) and 250 cosine similarity ( $cos\theta$ ) to capture both diversity and alignment in language use.

TTR is a measure of lexical diversity that compares the number of unique lemmas (types) to the total number of words (tokens) in a text ( $TTR = \frac{Types}{Tokens}$ ). This metric provides insight into the variety of vocabulary used in the responses, serving as an important indicator of linguistic complexity and richness. TTR ratios were computed within each participant's response as well as across students' responses to assess the lexical diversity of the verbal responses both within and between students.

257  $cos\theta$  measures the semantic similarity between words represented in vector embeddings. A pre-258 trained French model from FastText (<u>https://github.com/facebookresearch/fastText/</u>) was used 259 to provide a vectorial representation of word meanings. Cosine angle was then used to measure 260 the cosine similarity between words vectors =  $\frac{A.B}{\|A\|\|B\|}$ , offering a measure of how similar two words are in terms of their direction in the embedding space. We defined and calculated several specific similarity metrics, which were averaged across comparisons to generate a single similarity value per participant or odor:

- *Similarity Within Student*: This metric evaluates how similar the words used by
   a single participant were to each other, providing insights into the consistency
   of language within individual responses.
- 267 Similarity Between Students: This metric compares the vocabulary used by
   268 different students, to identify common themes and vocabulary.
- Similarity with an Expert Teacher: This metric assessed how closely student
   language aligned with the target language, by comparing student responses to a
   reference set of words representing an expert teacher's vocabulary,.
- *Similarity by Odor*: This metric groups words by the odors they described,
  analyzing how consistently students described the same odor.

274 Odor and visual recognition tasks: Recognition memory performance was assessed using 275 parameters derived from signal detection theory (Lockhart and Murdock, 1970). Four response 276 categories were defined based on the experimental conditions (target or distractor items) and 277 the students' behavioral responses (yes or no): Hit and Miss occurred when the target items 278 were correctly recognized or incorrectly rejected, respectively, and correct rejection (CR) and 279 false alarm (FA) occurred when the distractor items were correctly rejected or incorrectly 280 recognized, respectively. Two parameters were calculated from the Hit and FA scores: a 281 memory score (d'L) and a response bias score (CL). These scores were determined as follows:

282 
$$d'L = ln \left(\frac{HR(1 - FR)}{FR(1 - HR)}\right)$$

283 
$$CL = -0.5(ln\left(\frac{HR(1-HR)}{FR(1-FR)}\right))$$

Here, HR represents the Hit rate  $\left[\frac{Hit+0.5}{N1+1}\right]$ , FR represents the false alarm rate  $\left[\frac{FA+0.5}{N2+1}\right]$ , and N1 284 285 and N2 represent the number of target and distractor odorants, respectively, for which the 286 students provided a response. Memory scores can be positive or negative, indicating students' 287 good or poor performance in recognizing target items and rejecting distractor items, 288 respectively. A d'L score of 0 indicates chance-level performance, where students cannot 289 reliably distinguish between target and distractor items. Response bias scores indicate three 290 individual attitudes: conservative (tending to respond "no" to an item), neutral (responding 291 "yes" or "no" with equal probability), or liberal (tending to respond "yes"), with positive, 292 neutral, or negative values, respectively (Snodgrass and Corwin, 1988).

293 Categorization task and Categories verbalization: To analyze the categorization task, we 294 computed the number of groups created by students, as well as the agreement between cluster 295 assignments made by students and an expert teacher. To evaluate the consistency and reliability 296 of the clustering results across time points, two metrics were employed: Cohen's Kappa ( $\kappa$ ) and 297 Adjusted Rand Index (ARI). Cohen's Kappa measures the inter-rater agreement for categorical items, correcting for the agreement that could happen by chance. It is defined as  $\kappa = \frac{po-pe}{1-pe}$ 298 299 where  $p_o$  is the observed agreement and  $p_e$  is the expected agreement by chance. The Adjusted 300 Rand Index measures the similarity between two data clusterings, adjusting for the chance grouping of elements. It is defined as  $ARI = \frac{RI - E[RI]}{(RI) - E[RI]}$  where RI is the Rand Index, which 301 302 measures the percentage of decisions that are correct, and E[RI] is the expected value of the 303 Rand Index. These metrics were calculated for each participant by comparing their cluster 304 assignments with those of all other students and with the reference (expert teacher) at each time 305 point.

306 The descriptions of each category created by students were analyzed using part-of-speech 307 tagging, TTR scores and  $cos\theta$  metrics as in the evocation and description task (see the extensive 308 presentation of these scores above).

309 *Odor discrimination*: For each trio of odors, we determined whether the responses were correct 310 or incorrect. These responses were then summed for each participant at each time point to yield 311 a global score, referred to as the odor discrimination score. This score enabled to evaluate 312 overall odor discrimination performance over time, with a random level baseline set at 33%, 313 indicating chance-level performance.

314 Temporal dynamic of learning: In the final step of the analyses, we created the performance 315 profile of each participant at each time measurement by combining their scores in all metrics 316 significantly modulated by time. These profiles were normalized independently for each metric 317 to ensure meaningful comparisons between metrics. We then computed the Euclidean distances 318 between individual performance profiles across different time points (T0 vs. T1, T1 vs. T2, and 319 T0 vs. T2) to quantify the progress in performance over time. For example, the Euclidean 320 distance between participant profiles at time T0 and T1 containing three metrics (m1, m2, m3) 321 is given by the formula:

322 
$$d(T0,T1) = \sqrt{(m1_{T0} - m1_{T1})^2 + (m2_{T0} - m2_{T1})^2 + (m3_{T0} - m3_{T1})^2}$$

where  $m1_{T0}$ ,  $m2_{T0}$ ,  $m3_{T0}$  and  $m1_{T1}$ ,  $m2_{T1}$ ,  $m3_{T1}$  are the values of the three metrics at time points T0 and T1, respectively. Additionally, we calculated the Euclidean distance between each pair of students at each distinct time point (T0, T1, and T2) to assess the heterogeneity of profiles over time. For example, the Euclidean distance between the performance profiles containing three metrics (m1, m2, m3) of students J1 and J2 is given by the formula:

328 
$$d(J_{1},J_{2}) = \sqrt{(m_{J_{1}} - m_{J_{2}})^{2} + (m_{J_{1}} - m_{J_{2}})^{2} + (m_{J_{1}} - m_{J_{2}})^{2}}$$

#### 329 **2.7. Machine-learning analysis**

In the final step of our analysis, we aimed to predict students' training levels (T0: no training,
T1: 6 months of training, T2: 1.5 years of training) using students' performance profiles,
combining their scores on metrics significantly affected by time.

*Data Preprocessing:* Each metric was independently scaled by removing the mean and scaling
 to unit variance, and missing values were imputed using SimpleImputer (<u>https://scikit-</u>
 <u>learn.org/1.5/modules/generated/sklearn.impute.SimpleImputer.html</u>) with a mean strategy.
 This preprocessing ensured a normalized and complete dataset suitable for machine learning
 analysis.

338 Data Splitting and Cross-Validation: To maintain data integrity and prevent leakage, a nested
 339 cross-validation procedure was employed using a 10-fold Stratified Group K-Fold
 340 (StratifiedGroupKFold; <u>https://scikit-</u>

341 <u>learn.org/1.5/modules/generated/sklearn.model\_selection.StratifiedGroupKFold.html</u>) for both 342 the inner loop, which focused on hyperparameter tuning and model selection, and the outer 343 loop, which assesses the overall model performance, ensuring samples from the same individual 344 (grouped by 'Person ID') were not split across folds. This process was repeated 10 times to 345 obtain different random splits and provide a more robust estimate of the models' performance.

346 Models and Hyperparameter Tuning: We selected two machine learning models: Random 347 Forest (RF) and Linear Support Vector Machine (LSVM). RF was chosen for its ability to 348 handle complex feature interactions and provide valuable feature importance scores. LSVM is 349 effective for high-dimensional spaces, robust against overfitting, and particularly suitable for 350 small datasets, using a regularization parameter to manage noisy data. Hyperparameter tuning 351 performed GridSearchCV was using (https://scikit-352 learn.org/1.5/modules/generated/sklearn.model\_selection.GridSearchCV.html) with the same 353 StratifiedGroupKFold cross-validation strategy to find the best hyperparameters for each model. For RF, the hyperparameters tuned included n\_estimators, max\_depth,
min\_samples\_split, and min\_samples\_leaf. For LSVM, the hyperparameters tuned included C
and max\_iter.

357 Model Evaluation: Each tuned model was evaluated using 10-fold StratifiedGroupKFold 358 cross-validation on the training set to estimate generalization performance, with this process 359 repeated 10 times for robustness. The models were then trained on the entire training set and 360 tested on the hold-out test set, resulting in 100 total tests. Evaluation metrics, including accuracy 361 (the overall percentage of correct predictions) and F1-score (a harmonic mean of precision and 362 recall, where precision measures the proportion of correctly predicted instances of a class out 363 of all predicted instances, and recall measures the proportion of correctly predicted instances of 364 a class out of all actual instances of that class), were collected from the test sets of the outer 365 cross-validation loop, along with confusion matrices to visualize performance. These metrics 366 and the confusion matrix provide a comprehensive evaluation of the model's performance 367 across different classes, highlighting both strengths and areas for improvement. The models 368 were then trained on the entire training set and tested on the hold-out test set, resulting in 100 369 total tests. Evaluation metrics, including accuracy (the overall percentage of correct predictions) 370 and F1-score (a harmonic mean of precision and recall, where precision measures the proportion 371 of correctly predicted instances of a class out of all predicted instances, and recall measures the 372 proportion of correctly predicted instances of a class out of all actual instances of that class), 373 were collected from the test sets of the outer cross-validation loop, along with confusion 374 matrices to visualize performance.

375 *Feature Importance Analysis:* To understand the contribution of each feature to the models' 376 predictions, we extracted feature importances from the RF model and the absolute values of the 377 LSVM coefficients from the models trained during the outer cross-validation splits. Feature 378 importances were averaged, and standard deviations were computed to assess the models' variability. This analysis identifies the most significant features contributing to predictingtraining levels across both models and assesses the consistency of the results.

381 **2.8. Statistical data analysis** 

All reported metrics in Section 2.5 were computed independently for each measurement time 382 383 (T0, T1, and T2). We used repeated measures ANOVAs to compare these metrics across the 384 three time points, followed by Tukey post-hoc tests for significant results to identify specific 385 differences. One-sample t-tests were used to test whether performance was significantly 386 different from the chance level in recognition and discrimination tasks. In the machine-learning 387 analysis, one-tailed significance tests were performed to estimate the likelihood that the chance 388 level (33% for a 3-class problem) is the mean of the distribution of decoding accuracies from 389 multiple train/test splits. Additionally, a paired t-test was conducted to compare the accuracies 390 of the RF and LSVM models. On all reported-statistics, Bonferroni correction was applied to 391 all p-values for multiple comparisons. A p-value < 0.050 was considered as demonstrating 392 significant effect.

393

#### 2.9. Software and libraries for data analysis

394 All data preprocessing and analysis were performed using Python (https://www.python.org/). 395 The following packages and libraries were utilized: numpy (Harris et al., 2020) and pandas 396 (McKinney, 2010; The pandas development team, 2024) for data handling and manipulation, 397 SpaCy (Honnibal and Montani, 2017) for tokenization, Part-of-speech tagging, and 398 lemmatization for lemmatization, matplotlib (Hunter, 2007) and seaborn (Waskom, 2021) for 399 creating visualizations, and Scikit-learn for scaling, normalizing and ML analyses (Pedregosa 400 et al., 2011). These tools collectively facilitated comprehensive text processing, linguistic 401 analysis, and visualization. Jamovi software (The jamovi project, 2024), a free open-source 402 statistical tool with a user-friendly graphical interface built on R, was used for statistical403 analyses.

#### 404 **3. Results**

#### 405 **3.1. Description and Evocation tasks**

#### 406 *Description task*

The average number of words used in odor descriptions was assessed over three time periods (Figure 1A). We found a significant effect of time on the number of words used (F(2, 78) = 40.30, p < 0.001). Specifically, the average number of terms was lowest at T0 and progressively increased through T1 to T2. The number of words increased significantly from T0 to T1 (Mean Difference = -0.950, SE = 0.353, p = 0.031) and from T1 to T2 (Mean Difference = -1.713, SE = 0.242, p < 0.001), as well as from T0 to T2 (Mean Difference = -2.663, SE = 0.297, p < 413 0.001).

414 We then analyzed the average number of nouns, adjectives, and verbs in the odor descriptions 415 over time (Figures 1B & 1C). All descriptions consisted of word lists, with no actual 416 sentences present, resulting in the absence of any verbs. Similar to the trend observed with the 417 number of words, there was an effect of time on the number of nouns (F(2, 78) = 3.31, p =418 0.042), with a significant increase in the number of nouns from T0 to T2 (Mean Difference = 419 -1.225, SE = 0.476, p = 0.042), but no significant differences between T0 and T1 (Mean 420 Difference = -0.550, SE = 0.533, p = 0.562) and T1 and T2 (Mean Difference = -0.675, SE = 421 0.415, p = 0.247). The number of adjectives was also significantly influenced by time (F(2, 422 78) = 48.50, p < 0.001). An incremental increase in the use of adjectives over time was 423 demonstrated with an increase from T0 to T1 (Mean Difference = -1.48, SE = 0.381, p < 424 0.001) and from T1 to T2 (Mean Difference = -2.52, SE = 0.407, p < 0.001), as well as from



T0 to T2 (Mean Difference = -4.00, SE = 0.443, p < 0.001). In conclusion, students used more 426 words, nouns, and adjectives to describe odors as their learning progressed.

428 Figure 1. Odor description task. Increases in the average number of words (A), nouns (B), and adjectives (C) over 429 time (T0, T1 and T2) when describing odors. Error bars represent standard deviations. \* indicates p < 0.05, and 430 \*\*\* indicates p < 0.001.

431 In a second step, we explored the changes in lexical diversity within and between students' odor 432 descriptions over time using the TTR diversity score and  $cos\theta$  semantic similarity. While the 433 within-students TTR remained close to 1, indicating maximal diversity, and was relatively 434 stable over time (F(2, 78) = 1.60, p = 0.209; Figure 2A), within-participant semantic similarity significantly changed over time (F(2, 74) = 14.90, p < 0.001; Figure 2C). There were significant 435 436 increases in semantic similarity between T0 and T1 (Mean Difference = -0.081, SE = 0.030, p = 0.031), T1 and T2 (Mean Difference = -0.078, SE = 0.029, p = 0.031), and T0 and T2 (Mean 437 438 Difference = -0.160, SE = 0.029, p < 0.001). These results suggest that, over time, students used 439 a distinct yet increasingly specific and semantically similar set of words to describe each odor. 440 The between-student TTR diversity score changed over time (F(2, 78) = 5.23, p = 0.007; Figure

441 2B), with a significant decrease in lexical diversity between T0 and T2 (Mean Difference = 442 0.146, SE = 0.042, p = 0.004), but no significant differences between T0 and T1 or T1 and T2 (p > 0.245). We also found a significant effect of time on between-participant semantic 443

425

similarity (F(2, 78) = 90.90, p < 0.001), with incremental increases across the three time points: T0 and T1 (Mean Difference = -0.082, SE = 0.013, p < 0.001), T1 and T2 (Mean Difference = -0.065, SE = 0.010, p < 0.001), and T0 and T2 (Mean Difference = -0.147, SE = 0.010, p < .001) (Figure 2D). This result reflects an increase in the specificity of words used by students over time, confirming greater homogeneity in their vocabulary.

449 In addition, we examined the semantic similarity between students' answers and an expert 450 teacher and found a significant effect of time (F(2, 78) = 5.82, p = 0.004; Figure 2E), with 451 significant increase between T0 and T2 (Mean Difference = -0.073, SE = 0.022, p = 0.006), but 452 no significant differences between T0 and T1 (p = 0.191) and T1 and T2 (p = 0.418). Finally, 453 the analysis of the semantic similarity of descriptions for each odor across the three time points 454 demonstrated a significant effect of time (F(2, 22) = 21.50, p < 0.001; Figure 2F), with odors 455 being more similarly described at T2 than T0 (Mean Difference = -0.164, SE = 0.027, p < 0.001) and T1 (Mean Difference = -0.154, SE = 0.034, p = 0.002), but no significant difference 456 457 between T0 and T1 (Mean Difference = -0.010, SE = 0.022, p = 1.000). These results indicate 458 that over time, students' descriptions become more aligned with the expert's descriptions and 459 more consistent in how each odor is described, corroborating the increased precision and 460 standardization in their vocabulary.



Figure 2. Odor description task. Changes in lexical diversity using Type-Token Ratio (TTR) within (A) and between (B) students, and changes in semantic similarity within (C) and between (D) students, with an expert teacher (E), and by odor (F) over time (T0, T1, T2) when describing odors. Error bars represent standard deviations. \* indicates p < 0.05, \*\* indicates p < 0.01, and \*\*\* indicates p < 0.001.

### 466 Evocation task

461

467 No significant effect of time was found on the average number of words (Mean = 1.673, SE =

468 0.514, F(2, 78) = 2.06, p = 0.134), nouns (Mean = 1.228, SE = 0.452, F(2, 78) = 0.47, p = 0.47)

469 0.630), adjectives (Mean = 0.290, SE = 0.221, F(2, 78) = 0.83, p = 0.439), and verbs (Mean =

470 0.111, SE = 0.206, F(2, 78) = 1.43, p = 0.245) in the evocation responses.

471 The analysis of lexical diversity within-students' odor evocations using TTR diversity score did 472 not reveal significant differences of TTR over time (TTR<sub>within</sub>: Mean = 0.97, SE = 0.07, F(2, 473 (78) = 1.080, p = 0.346), but showed a significant effect of time on semantic similarity (F(2, 76)) 474 = 6.04, p = 0.004; Figure 3A), with similarity being higher at T2 than T0 (Mean Difference = -475 0.087, SE = 0.025, p = 0.004) but similar between T0 and T1 (Mean Difference = -0.031, SE = 476 0.024, p = 0.637) and T1 and T2 (Mean Difference = -0.056, SE = 0.026, p = 0.119). This result 477 suggests that as students progressed in their training, the words they used became increasingly 478 semantically similar.

479 When comparing between-participant lexical diversity, we found that while TTR diversity score 480 showed no effect of time (TTR<sub>between</sub>: Mean = 0.416, SE = 0.193, F(2, 78) = 1.080, p = 0.337), 481 semantic similarity significantly changed over time (F(2, 78) = 35.00, p < 0.001; Figure 3B), 482 with significant increases in similarity between T0 and T1 (Mean Difference = -0.026, SE = 483 0.009, p = 0.014), T1 and T2 (Mean Difference = -0.053, SE = 0.011, p < 0.001), and T0 and 484 T2 (Mean Difference = -0.079, SE = 0.009, p < .001). This result indicates increased similarity 485 in the students' vocabulary over time, reflecting greater semantic homogeneity in the words 486 they used.

487 In addition, the semantic similarity of evocations for each odor significantly evolved with time 488 (F(2, 22) = 11.20, p < 0.001; Figure 3C), with significant decreases between T0 and T1 (Mean 489 Difference = 0.063, SE = 0.015, p = 0.005) and T2 and T1 (Mean Difference = -0.105, SE = 490 0.023, p = 0.002), but no difference between T0 and T2 (Mean Difference = -0.042, SE = 0.027, 491 p = 0.462). These results indicate that over time, students' evocations were less consistent at 492 T1 than at T0 and T2, suggesting an effect of time unrelated to their level of training. Finally, 493 given that the evocation task is related to personal experiences and memories, we did not run 494 similarity comparisons to a reference.



496 Figure 3. Odor evocation task. Changes in cosine similarity ( $\cos\theta$ ) within (A) and between (B) students and by 497 odor (C) over time (T0, T1, T2) during odor-evoked evocations. Error bars represent standard deviations. \* 498 indicates p < 0.05, \*\* indicates p < 0.01, and \*\*\* indicates p < 0.001.

499

## **3.2.Discrimination task**

500 The Odor Discrimination Score was evaluated at three time points (T0, T1, T2). At T0, the mean score was 45.00% (SD = 27.79%), which was significantly above the chance level of 501 502 33.33% (t(39) = 2.66, p = 0.011). At T1, the mean score was 59.17% (SD = 25.58%), also significantly above chance (t(39) = 6.39, p < 0.001). At T2, the mean score was 53.33% (SD = 503 28.04%), and it remained significantly above the chance level (t(39) = 4.51, p < 0.001). 504

505 In addition to these comparisons against chance level, the within-students effect of time was 506 evaluated. The analysis revealed a marginally significant effect of time on performance (F(2,507 78) = 3.120, p = 0.050), with no significant differences between T0 and T1 (Mean Difference = -0.425, SE = 0.189, p = 0.091), T0 and T2 (Mean Difference = -0.250, SE = 0.159, p = 0.373), 508 and T1 and T2 (Mean Difference = 0.175, SE = 0.164, p = 0.874). 509 510 Overall, these results suggest that students' performance at each time point was significantly

511 better than random guessing. However, while Odor Discrimination Scores tend to fluctuate over

512 time, the differences between time points did not reach significance.

#### 513 **3.3. Olfactory and visual recognition tasks**

#### 514 *Odor recognition task*

515 The students were presented with target and distractor odors and were asked if they had smelled 516 them the day before. Memory scores were high (d'L: Mean = 1.669; SE = 1.385), indicating 517 proficiency in recognizing old odors and rejecting new ones at all measurement times (Figure 518 4A). d'L scores were significantly above chance at each measurement time: T0 (Mean = 1.324, 519 SD = 1.281, t(39) = 6.54, p < 0.001), T1 (Mean = 1.631, SD = 1.497, t(39) = 6.89, p < 0.001), 520 and T2 (Mean = 2.051, SD = 1.301, t(39) = 9.97, p < 0.001), demonstrating robust recognition 521 performance across time points. 522 However, the increase in d'L over time did not reach statistical significance (F(2,78) = 3.00, p 523 = 0.056). When examining correct response categories individually, a significant effect of time

was found for correct rejections (CR) (F(2,78) = 4.35, p = 0.016; Mean = 6.858; SD = 1.789;

Figure 4B), with a significant increase in the number of CRs between T0 and T2 (Mean

Difference = -1.00, SE = 0.36, p = 0.025). No effect of time was observed on the number of

Hits (F(2,78) = 0.78, p = 0.460; Mean = 6.875; SD = 1.771). The response bias was close to

zero at all time points, indicating no specific response strategy [ $CL_{T0}$ : Mean = -0.055, SD = 0.525;  $CL_{T1}$ : Mean = -0.028, SD = 0.518;  $CL_{T2}$ : Mean = 0.085, SD = 0.624] and was not

530 influenced by time (F(2,78) = 0.89, p = 0.428).



531

524

525

526

527

- 532 Figure 4. Recognition Performance in Olfactory and Visual Tasks. (A) Mean d'L scores for the olfactory task over
- 533 time. (B) Mean number of correct rejections (CR) for the olfactory task over time (10 distractors). (C) Mean d'L
- 534 scores for olfactory and visual tasks. Error bars represent standard deviations. '#' indicates a trend (p = 0.056), \*
- 535 indicates p < 0.05, and \*\*\* indicates p < 0.001.

536 Visual recognition task

- 537 In the visual recognition task, memory scores were very high, indicating strong proficiency in 538 recognizing old images and rejecting new ones at all measurement times (d'L: Mean = 4.068; 539 SD = 0.167). d'L scores were significantly above chance for each measurement time: T0 (Mean 540 = 3.789, SD = 1.898, t(39) = 12.63, p < 0.001), T1 (Mean = 4.403, SD = 1.429, t(39) = 19.49, 541 p < 0.001), and T2 (Mean = 4.011, SD = 2.083, t(39) = 12.18, p < 0.001), highlighting 542 consistently strong visual recognition performance.
- Visual memory scores were significantly higher than olfactory memory scores (F(1,39) = 131.54, p < 0.001; Figure 4C), illustrating that the visual memory task was easier for students than the olfactory memory task. However, for the visual task, there was no effect of time on either d'L memory scores or CL individual response categories (F's < 1.90; p's > 0.152). Finally, the response bias was close to zero at all time points (CL<sub>T0</sub>: Mean = -0.017, SD = 0.734; CL<sub>T1</sub>: Mean = -0.014, SD = 0.508; CL<sub>T2</sub>: Mean = -0.054, SD = 0.613), indicating no specific response strategy, and was not influenced by time (F(2,78) = 0.22, p = 0.807).

#### **3.4. Odor categorization task and Categories description**

#### 551 Odor categorization task

No significant effect of time was shown on the number of groups created by students (F(2, 76) = 1.65, p = 0.20; Mean = 6.342; SD = 1.738). The analysis of inter-student agreement indicated no agreement regarding the categories of odors created [Cohen's Kappa: Mean = 0.008, SD = 0.024] and showed no significant effect of time on agreement (F(2, 76) = 0.249, p = 0.780). However, the similarity of odors pairwise clustering (ARI) between students significantly varied with time (F(2, 76) = 17.30, p < 0.001; Figure 5A). There were significant increases between T0 and T1 (Mean Difference = -0.040, SE = 0.009, p < .001) and T0 and T2 (Mean Difference = -0.061, SE = 0.011, p < 0.001), but no significant difference between T1 and T2 (Mean Difference = -0.021, SE = 0.012, p = 0.227).

In addition, the analysis also revealed no agreement between students and the expert teacher 561 562 regarding the categories of odors created [Cohen's Kappa: Mean = -0.034, SD = 0.100], and no significant effect of time was found (F(2, 76) = 0.863, p = 0.426). However, the similarity of 563 564 odors pairwise clustering between students and the expert teacher significantly differed with 565 time (F(2, 76) = 3.260, p = 0.044), with a significant increase between T0 and T1 (Mean Difference = -0.066, SE = 0.022, p = 0.015), but no significant differences between T0 and T2 566 567 (Mean Difference = -0.052, SE = 0.029, p = 0.233) and between T1 and T2 (Mean Difference 568 = 0.014, SE = 0.030, p = 1.000).

569 Overall, while categorical agreement remained stable over time between students and the expert 570 teacher, the pairwise similarity of clustering increased with training, both among students and 571 between students and the expert teacher.



572

573 Figure 5. Odor categorization task. Changes of between-students (A) pairwise odors clustering similarity (ARI),

574 (B) lexical diversity (TTR), and (C) semantic similarity  $(cos(\theta))$  scores between students over time (T0, T1, T2).

575 Error bars represent standard deviations. \*\*\* indicates p < 0.001.

576 Description of odor categories

A significant effect of time on the average number of words was found (F(2,76) = 3.36, p = 0.040), showing a slight increase between T0 (Mean = 1.749, SD = 0.558) and T1 (Mean = 2.090, SD = 0.922) but no significant differences were found with T3 (Mean = 1.981, SD = 0.759;  $p'_s > 0.25$ ). No effect of time was revealed for the average number of nouns (Mean = 3.949, SD = 2.137; F(2,76) = 1.580, p = 0.213), nor adjectives (Mean = 6.017, SE = 2.626; F(2,76) = 0.514, p = 0.600). No verb was found in the verbal descriptions.

583 The between-participant lexical diversity (TTR) and semantic similarity  $(cos(\theta))$  significantly 584 differed over time (TTR<sub>between</sub>: F(2,76) = 14.800, p < .001;  $cos(\theta)_{between}$ : F(2,76) = 50.800, p < .001;  $cos(\theta)_{between}$ : F(2,76) = 50.800;  $cos(\theta)_{between}$ ; F(2,76) = 50.800;  $cos(\theta)_{between}$ ; F(2,76) = 50.800; F(585 .001; Figure 5B-C). There were significant decreases of lexical diversity between T0 and T1 586 (TTR Mean Difference = 0.20, SD = 0.04, p < 0.001) and T0 and T2 (TTR Mean Difference = 587 0.196, SD = 0.041, p < 0.001) but not between T1 and T2 (TTR Mean Difference = -0.002, SD = 0.039, p = 1.000). Conversely, we found significant increases of semantic similarity between 588 589 T0 and T1 (Mean Difference = -0.101, SD = 0.012, p < 0.001), T0 and T2 (Mean Difference = 590 -0.091, SD = 0.011, p < 0.001), but not between T1 and T2 (Mean Difference = 0.009, SD = 591 0.009, p = 0.984). Given that the principle of the categorization task is to create different groups 592 with specific descriptions, we did not run TTR and similarity comparisons within-students' 593 categories.

To sum-up, while the average number of words in category descriptions slightly increased over time, no significant changes were observed in the use of nouns or adjectives. However, there was a significant decrease in lexical diversity and a significant increase in semantic similarity between students' verbalizations over time, indicating a trend towards more homogeneous and semantically aligned verbal descriptions.

#### 599 **3.5.Temporal dynamic of learning**

600 Students' progress dynamics were analyzed over time. Students' performance profile was 601 created by combining their scores on the 14 metrics significantly affected by time (Table S2). 602 The within-students effect of time was found to be significant (F(2, 76) = 10.66, p < 0.001; 603 Figure 6A), with a linear and consistent increase in performance over time. Significant 604 differences in progress were observed between T0-T1 and T0-T2 (Mean Difference = -0.685, 605 SE = 0.187, p = 0.002), and between T1-T2 and T0-T2 (Mean Difference = -1.027, SE = 0.210, p < 0.001). However, no significant difference was found between T0-T1 and T1-T2 (Mean 606 607 Difference = 0.341, SE = 0.254, p = 0.515). Overall, these results suggest that students' progress 608 was consistent from T0 to T1 and from T1 to T2.



#### 609

610 Figure 6. Temporal dynamics of students' learning. (A) Changes of within-students' performance over reflecting 611 their progress as their training progressed, and (B) Changes of between-students' performance over time showing 612 the heterogeneity of performance between students at each time (T0, T1, T2). Error bars represent standard 613 deviations. \*\*, indicates p < 0.01; \*\*\* indicates p < 0.001

The heterogeneity of performance between students was evaluated based on the score differences between each pair of students at each time measurement (T0, T1, and T2). The effect of time on the heterogeneity of performance was significant (F(2, 1480) = 29.00, p < 617 0.001; Figure 6B), with a greater heterogeneity at T1 compared to T0 (Mean Difference = 0.241, 618 SE = 0.064, p < 0.001) and T2 (Mean Difference = 0.429, SE = 0.054, p < 0.001). However, 619 heterogeneity was lower at T2 than at T0 (Mean Difference = -0.188, SE = 0.050, p < 0.001). 620 These results suggest that the initial heterogeneity of students' profiles at T0 increased by T1, 621 indicating varying speeds of learning among students. By T2, the heterogeneity decreased to a 622 level even below the initial point, reflecting a convergence in the learning levels of students.

623

#### **3.6.** Predicting students' training levels

624 Given the clear changes in performance over time, we explored whether it was possible to 625 identify students' level of training with high accuracy based on their performance across tasks 626 using Random Forest (RF) and Linear Support Vector Machine (LSVM) classifiers.

627 On average, the best performance was obtained with the RF classifier, which correctly 628 identified the students' training level 78.9% of the time, with a margin of error of  $\pm 11.1\%$ 629 (Figure 7A). The statistical analysis showed that the accuracy score was highly significant (t = 630 41.30; one-tailed corrected p-value < 0.001), and the model achieved an F1-score of 0.783  $\pm$ 631 0.113, indicating strong performance in both precision (accuracy of positive predictions) and 632 recall (ability to find all relevant instances). For the RF classifier, the most crucial features were 633 the semantic similarities of odor descriptions among students (description sim-BTW) and the 634 semantic similarities verbalizations between odor category among students 635 (categorisation sim-BTW; Figure 7B).

The SVM classifier also performed well, correctly identifying the training level of students 74.2% of the time, with a margin of error of  $\pm$  6.9% (t = 59.173; one-tailed corrected p-value < 0.001; Figure 7A), and an F1-score of 0.728  $\pm$  0.076, indicating good overall performance of the model across all classes. The feature importance analysis for the SVM classifier (Figure 7C) also highlighted the importance of the semantic similarities between students' responses 641 measured in the description and the verbalization part of the categorization tasks in determining



642 the training level of students.



648 Rec-olf: odor recognition task.

643

649 When directly comparing the performance of the RF and LSVM classifiers, the model 650 comparison indicated that the RF model significantly outperformed the LSVM model (t = 651 4.101; p < 0.001). This suggests that the RF model, with its large number of estimators and 652 fully grown trees, was better suited for capturing the underlying patterns in the perceptual and 653 cognitive data, leading to its superior performance compared to the LSVM classifier. The full 654 hyperparameters tuning results are summarized in Table S3.

655 **4. Discussion** 

The primary objective of this research was to examine how academic olfactory training during a 1.5-year educational program impacts the development of olfactory expertise in perfumery students. The study focused on a broad range of tasks, including odor description, evocation, recognition, discrimination, and categorization. Our main findings reveal that while students demonstrated significant improvement in their ability to describe and categorize odorsreflected in richer (words, nouns, adjectives count) and more consistent language (TTR scores and  $\cos\theta$  similarity), and its greater alignment with expert terminology ( $\cos\theta$  similarity)—their non-verbally mediated abilities, such as odor discrimination and recognition accuracies, did not show substantial enhancement over time.

665 Our results demonstrate that as students progressed through the training program, their 666 vocabulary and descriptive abilities expanded, evidenced by an increased use of words, nouns, 667 and adjectives to describe odors. Over time, however, lexical diversity within the student group 668 decreased, while semantic similarity increased-both among students and between students and 669 the expert teacher. This indicates a homogenization of vocabulary, as students adopted a more 670 standardized language for odor description. This convergence suggests that the training not only 671 enriched their descriptive repertoire but also guided them toward a shared, expert-aligned 672 terminology. Interestingly, the number of words used in odor-evoked evocations did not change 673 significantly over time, yet semantic similarity increased within and between students. This 674 implies that while students did not elaborate more in terms of word count when recalling odor-675 associated memories, their language became more semantically aligned. This reflects the 676 influence of training on the conceptualization and verbal encoding of olfactory experiences, 677 facilitating more efficient retrieval and communication of olfactory information (Larsson and 678 Willander, 2009). These findings align with previous research indicating that professional 679 perfumers use more detailed and precise language focusing on chemical and olfactory qualities 680 (Sezille et al., 2014) and that olfactory cognitive abilities can be improved with training (Royet 681 et al., 2013).

In addition to improvements in descriptive abilities, our findings reveal significant changes in how students categorized odors over time. While the number of groups created and the categorical agreement among students did not change significantly, the similarity of pairwise odor clustering between students increased. This suggests that as students progressed through the training, they began grouping odors in increasingly similar ways, reflecting a convergence in their perceptual representation of odors. Moreover, although the average number of words used in category descriptions increased slightly, there was a noticeable decrease in lexical diversity and a significant increase in semantic similarity among students' verbalizations. In other words, students used fewer unique terms and described odors in ways that were more aligned with one another. The increased alignment in odor grouping and the standardization of vocabulary suggest that students are developing a shared terminology for categorizing odors.

693 Our results highlight the central role of language in structuring olfactory knowledge, and align 694 with prior cross-cultural findings highlighting the variability of odor lexicons and the 695 importance of consistent linguistic frameworks for olfactory description (Majid, 2021). Studies 696 on languages such as Jahai and Maniq demonstrate that dedicated lexical fields for odors can 697 emerge when olfactory experience is consistently emphasized and linguistically encoded 698 (Majid and Burenhult, 2014). These findings challenge the notion proposed by Olofsson and 699 Gottfried that olfactory naming difficulties are primarily rooted in neural anatomical constraints 700 (Olofsson and Gottfried, 2015). Instead, our results support the idea that developing olfactory 701 expertise and a shared vocabulary for odors is not only possible but also heavily influenced by 702 academic training practices.

703 Despite the observed improvements in verbally-mediated and categorization tasks, our study 704 did not reveal significant enhancements in non-verbally mediated abilities such as odor 705 discrimination and recognition. The odor discrimination scores showed no significant 706 improvement over time, and while there was a trend toward better odor recognition 707 performance, it did not reach statistical significance. These findings are consistent with research 708 by Filiz et al. (2022), who found no substantial improvement in discrimination or recognition 709 abilities in wine students over a similar period of training. This lack of improvement could be 710 attributed to several factors. First, high-order non-verbally mediated tasks such as odor

711 discrimination and recognition may require extensive and more prolonged exposure to a wider 712 variety of odors to see measurable changes (Rovet et al., 2013), which might not have been 713 fully achieved within the duration of our study. Second, prior research suggests that olfactory 714 training effects are often task-specific and may not generalize across different odor categories 715 (Bende and Nordin, 1997; Parr et al., 2002; Croijmans and Majid, 2016; Poupon et al., 2018). 716 At ISIPCA, students are trained using a structured curriculum that focuses on the gradual 717 introduction of raw materials classified by olfactory families, progressing from simpler to more 718 complex odors. While this approach strengthens categorical organization and verbal 719 descriptors, it may also reinforce expertise within specific odor categories rather than promoting 720 broader generalization across categories. For instance, training sessions often emphasize 721 identifying and describing fragrances or flavors within specific families, such as fruity, floral, 722 or toasted, which could limit the transfer of perceptual skills to novel or less familiar odor 723 groups. Another plausible factor is that the methods employed in this study, including task 724 design and stimuli selection, may not have been sensitive enough to capture subtle changes in 725 olfactory performance.

726 The use of natural language processing (NLP) and machine learning (ML) tools in this study 727 provided valuable insights into olfactory expertise, reinforcing the centrality of language in its 728 acquisition. These computational methods allowed us to capture subtle patterns in language use, 729 uncovering links between linguistic descriptors and expertise that would have been overlooked 730 with more traditional approaches. Features related to verbal tasks-such as the richness and 731 precision of odor descriptions, the use of specific olfactory terminology, and semantic similarity 732 measures in the description and categorization tasks-were the strongest predictors of students' 733 performance and expertise levels. Without NLP and ML techniques, the effects of olfactory 734 training on expertise would have gone undetected. This highlights the potential and importance

of NLP and ML methods for better characterizing olfactory language and expertise, as well as
refining training methodologies.

737 Our effort to better characterize the semantic properties of olfactory language aligns with other 738 computational approaches that have recently emerged to deepen our understanding of olfactory 739 perception and its connection to language and expertise. For instance, Iatrapoulos et al. (2018) 740 introduced new metrics-the Olfactory Association Index (OAI) and Olfactory Specificity 741 Index (OSI)-to quantitatively assess how strongly words are associated with olfaction and 742 how specifically they describe odors. Similarly, Hörberg et al. (2022) mapped the olfactory 743 vocabulary of English using NLP techniques, identifying semantic dimensions such as valence, 744 concreteness, and edibility that organize the olfactory lexicon. Another study by Hörberg et al. 745 (2025) compared sensory vocabularies across domains such as wine, perfume, and food using 746 computational methods, revealing critical differences in domain specificity, descriptor 747 preferences, and semantic dimensions. These studies highlight how computational analyses of 748 large-scale textual data can uncover the linguistic structures underlying sensory descriptions 749 and expertise.

Collectively, these computational approaches to olfactory language underscore the usefulness of NLP and ML tools not only for characterizing olfactory expertise but also for advancing training methodologies in olfactory and sensory sciences. By leveraging these innovative methods, we can gain deeper insights into the verbal underpinnings of olfactory expertise, and develop more effective strategies for training and assessment.

The olfactory training program at ISIPCA demonstrates considerable effectiveness in enabling students to develop and utilize a precise, shared language for describing and categorizing odors. Through structured exercises, systematic exposure, and standardized verbalization practices, students progressively refine their vocabulary, achieving greater alignment with expert terminology. While these findings underscore the program's proficiency in fostering linguistic and conceptual frameworks critical for olfactory expertise, there remains scope for further
refinement. Building upon this solid foundation, our results, supported by prior research,
suggest three key areas that could further optimize training outcomes.

763 Firstly, emphasizing vocabulary specificity and encouraging the consistent use of standardized 764 terminology can strengthen associative links between odors and their descriptions, facilitating 765 memory and conceptualization. The structured language training followed at ISIPCA, where 766 students repeatedly smelled, described, and classified raw materials using Jaubert's field of 767 odors (1995), already demonstrated how systematic exposure and verbalization promote the 768 development of a precise and shared olfactory lexicon. To further enhance this approach, the 769 introduction of a shared textual vocabulary dictionary could ensure that all students have access 770 to a unified reference for describing odors. This resource would serve as a foundational tool for 771 communication and conceptual alignment, reinforcing the consistent use of shared language 772 across both students and teachers.

773 Secondly, integrating mental imagery training into the curriculum could enhance cognitive 774 skills crucial for perfumers. Research by Plailly et al. (2012) showed that olfactory mental 775 imagery is central to perfumers' daily tasks and improves significantly with training. Croijmans 776 et al. (2020) demonstrated that mental olfactory imagery can be developed through targeted 777 practice rather than being an inherent ability. Similarly, Stevenson et al. (2007) found that 778 training in odor naming enhances mental imagery by providing verbal anchors that strengthen 779 cognitive associations. Given the structured exercises in ISIPCA training program, such as 780 identifying raw materials blindly and describing their volatility, incorporating explicit mental 781 imagery tasks could complement these activities and support students in imagining, recalling, 782 and categorizing odors more effectively. Testing these processes in future studies would clarify 783 how mental imagery integrates with verbal and perceptual learning to build olfactory expertise.

784 Thirdly, since improvements in perceptual abilities such as odor discrimination may be slower, 785 training programs could incorporate prolonged exposure to a wider variety of odors and more 786 intensive discrimination tasks. This could include differentiating subtle variations in odor 787 mixtures or thresholds, as suggested by Morquecho-Campos et al. (2019) and Poupon et al. 788 (2018). For instance, ISIPCA program already introduced students to a progression of odor 789 intensities and complex combinations through structured sessions, but future designs might 790 integrate targeted exercises to refine sensory acuity further. These could involve evaluating 791 persistence and volatility or focusing more extensively on odor combinations to sharpen 792 discrimination and perceptual organization.

Moreover, our study highlights the dynamics of learning over time and individual differences among students. Greater variability in performance observed at the midpoint of training (T1) suggests that individuals progress at different rates. By the end of the program (T2), performance levels converged, indicating that prolonged training promotes a more uniform level of expertise. These findings underscore the need for personalized training approaches that accommodate individual learning trajectories, tailoring intensity and content to optimize outcomes and help students reach their full potential.

### 800 Conclusion

801 In conclusion, our study highlights the pivotal role of language in the development of olfactory 802 expertise among perfumery students. Over the course of the 1.5-year training program, students 803 demonstrated significant improvements in their ability to describe and categorize odors, 804 developing a richer and more precise olfactory lexicon. These linguistic advancements were 805 accompanied by increased semantic alignment among students and with expert standards, 806 underscoring the effectiveness of structured training in fostering a shared vocabulary essential for professional expertise. However, the findings reveal limitations in the enhancement of non-807 808 verbally mediated abilities, such as odor discrimination and recognition. This disparity suggests 809 that while acquiring a shared olfactory language is crucial for developing conceptual and 810 descriptive skills, perceptual expertise may require more extensive or targeted sensory 811 exercises. Our study highlights the need for complementary strategies, including prolonged 812 exposure to diverse odor profiles, mental imagery training, and personalized training 813 approaches tailored to individual learning trajectories, to fully support the multifaceted nature 814 of olfactory expertise. Finally, our results also underscore the value of computational 815 approaches, such as natural language processing (NLP) and machine learning (ML), in 816 characterizing olfactory expertise and refining training methodologies. Future research should 817 explore these approaches to deepen our understanding of how olfactory expertise emerges.

818

### 5. Conflict of interests

819 The authors declare no conflicts of interest related to this research.

#### 820 **6.** Funding

The experimental costs associated with this research were fully covered by ISIPCA. The work of the researchers was supported by funding from their respective institutions. No additional external funding was received for this study.

#### 824 **7.** Acknowledgments

We would like to express our sincere gratitude to all the perfumery students who participated in this study for their time and effort. We also extend our heartfelt thanks to our colleagues at ISIPCA for their invaluable support in facilitating the data acquisition process. Their assistance and collaboration were instrumental in the successful completion of this research.

#### 829 8. Data availability

830 The data that support the findings of this study will be made publicly available in an831 anonymized format upon publication in the Open Science Framework (OSF).

- Al Aïn, S., Poupon, D., Hétu, S., Mercier, N., Steffener, J., and Frasnelli, J. 2019. Smell training
  improves olfactory function and alters brain structure. NeuroImage. 189:45–54.
- 835 Banks, S.J., Sreenivasan, K.R., Weintraub, D.M., Baldock, D., Noback, M., Pierce, M.E.,
- 836 Frasnelli, J., James, J., Beall, E., Zhuang, X., et al. 2016. Structural and Functional MRI
- Bifferences in Master Sommeliers: A Pilot Study on Expertise in the Brain. Front HumNeurosci. 10:414.
- 839 Barbeau, E., Didic, M., Tramoni, E., Felician, O., Joubert, S., Sontheimer, A., Ceccaldi, M.,
- and Poncet, M. 2004. Evaluation of visual recognition memory in MCI patients. Neurology.
  62:1317–1322.
- Bende, M., and Nordin, S. 1997. Perceptual Learning in OlfactionProfessional Wine Tasters
  versus Controls. Physiol Behav. 62:1065–1070.
- 844 Bensafi, M., Fournel, A., Joussain, P., Poncelet, J., Przybylski, L., Rouby, C., and Tillmann, B.
- 845 2017. Expertise shapes domain-specific functional cerebral asymmetry during mental imagery:
- the case of culinary arts and music. Eur J Neurosci. 45:1524–1537.
- 847 Carreiras, M., Quiñones, I., Chen, H.A., Vázquez-Araujo, L., Small, D., and Frost, R. 2024.
  848 Sniffing out meaning: Chemosensory and semantic neural network changes in sommeliers.
  849 Hum Brain Mapp. 45:e26564.
- 850 Castriota-Scanderbeg, A., Hagberg, G.E., Cerasa, A., Committeri, G., Galati, G., Patria, F.,
- 851 Pitzalis, S., Caltagirone, C., and Frackowiak, R. 2005. The appreciation of wine by sommeliers:
- a functional magnetic resonance study of sensory integration. NeuroImage. 25:570–578.

- 853 Croijmans, I., Arshamian, A., Speed, L.J., and Majid, A. 2021. Wine experts' recognition of
  854 wine odors is not verbally mediated. J Exp Psychol Gen. 150:545–559.
- 855 Croijmans, I., and Majid, A. 2016. Not All Flavor Expertise Is Equal: The Language of Wine
  856 and Coffee Experts. PLOS ONE. 11:e0155845.
- 857 Croijmans, I., Pellegrino, R., and Janice Wang, Q. 2024. Demystifying wine expertise through
- the lens of imagination: Descriptions and imagery vividness across sensory modalities. FoodRes Int Ott Ont. 182:114159.
- 860 Croijmans, I., Speed, L.J., Arshamian, A., and Majid, A. 2020. Expertise Shapes Multimodal
- 861 Imagery for Wine. Cogn Sci. 44:e12842.
- Belon-Martin, C., Plailly, J., Fonlupt, P., Veyrac, A., and Royet, J.P. 2013. Perfumers' expertise
  induces structural reorganization in olfactory brain regions. NeuroImage. 68:55–62.
- 864 Filiz, G., Poupon, D., Banks, S., Fernandez, P., and Frasnelli, J. 2022. Olfactory bulb volume
- and cortical thickness evolve during sommelier training. Hum Brain Mapp. 43:2621–2633.
- 866 Fournel, A., Sezille, C., Licon, C.C., Sinding, C., Gerber, J., Ferdenzi, C., Hummel, T., and
- Bensafi, M. 2017. Learning to name smells increases activity in heteromodal semantic areas.
  Hum Brain Mapp. 38:5958–5969.
- Haehner, A., Tosch, C., Wolz, M., Klingelhoefer, L., Fauser, M., Storch, A., Reichmann, H.,
  and Hummel, T. 2013. Olfactory Training in Patients with Parkinson's Disease. PLOS ONE.
  871 8:e61680.
- 872 Harris, C.R., Millman, K.J., Van Der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D.,
- 873 Wieser, E., Taylor, J., Berg, S., Smith, N.J., et al. 2020. Array programming with NumPy.
- 874 Nature. 585:357–362.

- 875 Herz, R. 2005. The unique interaction between language and olfactory perception and876 cognition.
- Honnibal, M., and Montani, I. 2017. spaCy 2: Natural language understanding with Bloom
  embeddings, convolutional neural networks and incremental parsing.
- 879 Hörberg, T., Kurfalı, M., and Olofsson, J.K. 2025. Chemosensory vocabulary in wine, perfume
- and food product reviews: Insights from language modeling. Food Qual Prefer. 124:105357.
- 881 Hörberg, T., Larsson, M., and Olofsson, J.K. 2022. The Semantic Organization of the English
- 882 Odor Vocabulary. Cogn Sci. 46:e13205.
- Hughson, A.L., and Boakes, R.A. 2001. Perceptual and cognitive aspects of wine expertise.
  Aust J Psychol. 53:103–108.
- 885 Hunter, J.D. 2007. Matplotlib: A 2D Graphics Environment. Comput Sci Eng. 9:90–95.
- 886 Iatropoulos, G., Herman, P., Lansner, A., Karlgren, J., Larsson, M., and Olofsson, J.K. 2018.
- The language of smell: Connecting linguistic and psychophysical properties of odor descriptors.
  Cognition. 178:37–49.
- Jaubert, J.-N., Tapiero, C., and Dore, J.-C. 1995. The Field of Odors: Toward a Universal
  Language for Odor Relationships. Perfum Flavorist. 20.
- 891 Larsson, M., and Willander, J. 2009. Autobiographical odor memory. Ann N Acad Sci.
  892 1170:318–23.
- Lesschaeve, I., and Issanchou, S. 1996. Effects of panel experience on olfactory memory
  performance: influence of stimuli familiarity and labeling ability of subjects. Chem Senses.
  21:699–709.

- Lockhart, R., and Murdock, B. 1970. Memory and the theory of signal detection. Psychol Bull.
  74:100–9.
- Majid, A. 2021. Human Olfaction at the Intersection of Language, Culture, and Biology. Trends
  Cogn Sci. 25:111–123.
- 900 Majid, A., and Burenhult, N. 2014. Odors are expressible in language, as long as you speak the
- 901 right language. Cognition. 130:266–270.
- 902 McKinney, W. 2010. Data Structures for Statistical Computing in Python. In: Austin, Texas.903 pp. 56–61.
- 904 Morquecho-Campos, P., Larsson, M., Boesveldt, S., and Olofsson, J.K. 2019. Achieving
- 905 Olfactory Expertise: Training for Transfer in Odor Identification. Chem Senses. 44:197–203.
- Olofsson, J.K., and Gottfried, J.A. 2015. The muted sense: neurocognitive limitations ofolfactory language. Trends Cogn Sci. 19:314–321.
- 908 Parr, W.V. 2019. Demystifying wine tasting: Cognitive psychology's contribution. Food Res
  909 Int. 124:230–233.
- Parr, W.V., Heatherbell, D., and White, K.G. 2002. Demystifying wine expertise: olfactory
  threshold, perceptual skill and semantic memory in expert and novice wine judges. Chem
  Senses. 27:747–755.
- 913 Pazart, L., Comte, A., Magnin, E., Millot, J.-L., and Moulin, T. 2014. An fMRI study on the
- 914 influence of sommeliers' expertise on the integration of flavor. Front Behav Neurosci. 8:358.

- 915 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
- 916 Prettenhofer, P., Weiss, R., Dubourg, V., et al. 2011. Scikit-learn: Machine Learning in Python.
- 917 J Mach Learn Res. 12:2825–2830.
- Plailly, J., Delon-Martin, C., and Royet, J.P. 2012. Experience induces functional
  reorganization in brain regions involved in odor imagery in perfumers. Hum Brain Mapp.
  33:224–234.
- Poupon, D., Fernandez, P., Archambault Boisvert, S., Migneault-Bouchard, C., and Frasnelli,
  J. 2018. Can the Identification of Odorants Within a Mixture Be Trained? Chem Senses.
  43:721–726.
- Royet, J.P., Plailly, J., Saive, A.L., Veyrac, A., and Delon-Martin, C. 2013. The impact of
  expertise in olfaction. Front Psychol. 4.
- Sezille, C., Fournel, A., Rouby, C., Rinck, F., and Bensafi, M. 2014. Hedonic appreciation and
  verbal description of pleasant and unpleasant odors in untrained, trainee cooks, flavorists, and
  perfumers. Front Psychol. 5:12.
- Snodgrass, J.G., and Corwin, J. 1988. Pragmatics of measuring recognition memory:
  applications to dementia and amnesia. J Exp Psychol Gen. 117:34–50.
- 931 Solomon, G.E.A. 1997. Conceptual Change and Wine Expertise. J Learn Sci. 6:41–60.
- 932 Sreenivasan, K., Zhuang, X., Banks, S.J., Mishra, V., Yang, Z., Deshpande, G., and Cordes, D.
- 933 2017. Olfactory Network Differences in Master Sommeliers: Connectivity Analysis Using
- 934 Granger Causality and Graph Theoretical Approach. Brain Connect. 7:123–136.
- Stevenson, R.J., Case, T.I., and Mahmut, M. 2007. Difficulty in evoking odor images: the role
  of odor naming. Mem Cognit. 35:578–589.

- 937 The jamovi project. 2024. Jamovi (Version 2.5) [Computer Software].
- 938 The pandas development team. 2024. pandas-dev/pandas: Pandas. Zenodo.
- 939 Vanek, N., Sóskuthy, M., and Majid, A. 2021. Consistent verbal labels promote odor category
- 940 learning. Cognition. 206:104485.
- 941 Waskom, M.L. 2021. seaborn: statistical data visualization. J Open Source Softw. 6:3021.
- 942 Zambom-Ferraresi, F., Zambom-Ferraresi, F., Fernández-Irigoyen, J., Lachén-Montes, M.,
- 943 Cartas-Cejudo, P., Lasarte, J.J., Casares, N., Fernández, S., Cedeño-Veloz, B.A., Maraví-
- Aznar, E., et al. 2021. Olfactory Characterization and Training in Older Adults: Protocol Study.
- 945 Front Aging Neurosci. 13.

946

#### 947 **10. Figure Legends**

Figure 1. Odor description task. Increases in the average number of words (A), nouns (B), and adjectives (C) over time (T0, T1 and T2) when describing odors. Error bars represent standard deviations. \* indicates p < 0.05, and \*\*\* indicates p < 0.001.

Figure 2. Odor description task. Changes in lexical diversity using Type-Token Ratio (TTR)
within (A) and between (B) students, and changes in semantic similarity within (C) and between
(D) students, with an expert teacher (E), and by odor (F) over time (T0, T1, T2) when describing
odors. Error bars represent standard deviations. \* indicates p < 0.05, \*\* indicates p < 0.01, and</li>

955 \*\*\* indicates p < 0.001.

**Figure3. Odor evocation task.** Changes in cosine similarity  $(\cos\theta)$  within (A) and between (B) students and by odor (C) over time (T0, T1, T2) during odor-evoked evocations. Error bars represent standard deviations. \* indicates p < 0.05, \*\* indicates p < 0.01, and \*\*\* indicates p < 0.001.

Figure4. Recognition Performance in Olfactory and Visual Tasks. (A) Mean d'L scores for the olfactory task over time. (B) Mean number of correct rejections (CR) for the olfactory task over time (10 distractors). (C) Mean d'L scores for olfactory and visual tasks. Error bars represent standard deviations. '#' indicates a trend (p = 0.056), \* indicates p < 0.05, and \*\*\* indicates p < 0.001.

965 **Figure5. Odor categorization task.** Changes of between-students (A) pairwise odors 966 clustering similarity (ARI), (B) lexical diversity (TTR), and (C) semantic similarity ( $cos(\theta)$ ) 967 scores between students over time (T0, T1, T2). Error bars represent standard deviations. \*\*\* 968 indicates p < 0.001.

969 Figure6. Temporal dynamics of students' learning. (A) Changes of within-students'
970 performance over reflecting their progress as their training progressed, and (B) Changes of

between-students' performance over time showing the heterogeneity of performance between
students at each time (T0, T1, T2). Error bars represent standard deviations. \*\*, indicates p <</li>
0.01; \*\*\* indicates p < 0.001</li>

Figure7. Feature Importances and Model Accuracies of students' level of training
predictions. (A) Test accuracies for the RF and LSVM models. (B) Top 6 feature importances
for the Random Forest (RF) classifier. (C) Top 6 feature importances for the Linear Support
Vector Machine (LSVM) classifier. Error bars represent standard deviations. \*\*\* indicates p <</li>
0.001. Adj: adjectives; Cat: categorization task; Desc: description task; NB: number; Rec-olf:
odor recognition task.

980

# 981 **11. Supplementary materials**

- 982 **Table S1.** Exhaustive list of olfactory stimuli
- 983 Note. Abs, absolute ; Ess, essence ; EO, essential oil. Perfumes are in italics. Hyperscript letters
- 984 referred to odorants that are used together in a mixture

Tasks	Odorants	Reference		
Evocation and	Vegetable	Formulated at the		
description*	Spicy grilled	ISIPCA using in-		
	Pharmaceutical	house recipes		
	Floral			
	Sun-soaked Sea scent			
	Detergent			
	Dried fruit honey			
	Musky animal scent			
	Late evening			
	Green tea			
	Fig			
	Bin scent			
Discrimination	Flowerbomb	Viktor & Rolf		
	La Vie Est Belle	Lancôme		
	La Nuit Trésor	Lancôme		
	Black Opium	Yves Saint Laurent Yves Saint Laurent		
	L'Homme			
	Pink Grape Fruit	Hermès		
Recognition T0	Ambrarome	Nactis		
	Thymol	Sigma Aldrich		
	Borneol	Prodasynth		
	Nutmeg	Robertet		
	Light blue	Dolce Gabbana		

	Le REM	Reminiscence
	Accord EMMA	ISIPCA Formulation
	Aromatics in black	Clinique
	Aromatics in white	Clinique
	L'Instant Magic	Guerlain
	Eau d'Hadrien	A.Goutal
	Apogée	Vuitton
	Contre moi	Vuitton
	Dans la peau	Vuitton
	Matière noire	Vuitton
	Mille feux	Vuitton
	Rose des vents	Vuitton
	Turbulences	Vuitton
	Piment brûlant	Artisan Parfumeur
	Coach	New York
Recognition T1	Dew fruit <sup>a</sup>	Givaudan
	Limon ess <sup>a</sup>	Albert Vieille
	Phenylacetic acid <sup>b</sup>	Sigma Aldrich
	beta-naphthyl methyl ketone <sup>b</sup>	BLH PIM
	Chocovan <sup>c</sup>	GIVAUDAN
	Nutmeg ess <sup>c</sup>	Albert Vieille
	Damascenia <sup>d</sup>	FIRMENICH
	Nootkatone <sup>d</sup>	PRODASYNTH
	Exaltolide <sup>d</sup>	FIRMENICH

Vertofix <sup>e</sup>	IFF
Alpha-Ionone <sup>e</sup>	Sigma Aldrich
Cashmeran <sup>f</sup>	IFF
Ethyl linalool <sup>f</sup>	Givaudan
Ginger ess <sup>f</sup>	Albert Vieille
Russian leather <sup>g</sup>	NACTIS
beta-naphthyl methyl ketone <sup>g</sup>	BLH PIM
Nerol <sup>h</sup>	IFF
Anisyl Acetate <sup>h</sup>	Givaudan
Manzanate <sup>i</sup>	Givaudan
Habanolide <sup>i</sup>	FIRMENICH
Magnolia flowers ess oil <sup>j</sup>	LMR
Cedramber <sup>j</sup>	IFF
Fruit sec $\mathbb{C}^k$	FIRMENICH
Adoxal <sup>k</sup>	Givaudan
Osmanthus abs <sup>1</sup>	BLH PIM
Tamarine base <sup>1</sup>	FIRMENICH
Verdox <sup>m</sup>	IFF
Cashmeran <sup>m</sup>	IFF
Veloutone <sup>n</sup>	FIRMENICH
Viridine <sup>n</sup>	BLH PIM
Sandalore <sup>o</sup>	Givaudan
Iso E Super <sup>o</sup>	IFF
Acetate cis 3 hexenyle <sup>p</sup>	Givaudan

	Rhodinol <sup>p</sup>	Givaudan
	Terebentine ess <sup>q</sup>	FIRMENICH
	Vanillin <sup>q</sup>	SOLVAY
	Scentenal <sup>r</sup>	FIRMENICH
	Miel Blanc © <sup>r</sup>	PCW
	Bergamot ess <sup>r</sup>	Payan & Bertrand
	Gamma Nonalactone <sup>s</sup>	Sigma Aldrich
	Allyl amyl glycolate <sup>s</sup>	IFF
	Cyclotene <sup>t</sup>	Sigma Aldrich
	Cuminic aldehyde <sup>t</sup>	Givaudan
Recognition T2	Lime EO <sup>u</sup>	Payan & Bertrand
	Bitter orange EO <sup>u</sup>	Robertet
	Verdox <sup>u</sup>	IFF
	Menthanyl acetate <sup>v</sup>	FIRMENICH
	Vertofix <sup>v</sup>	IFF
	Decanal <sup>v</sup>	Givaudan
	Linalyl acetate <sup>w</sup>	Sigma Aldrich
	Opoponax EO <sup>w</sup>	FIRMENICH
	Cis-3-Hexenol <sup>x</sup>	PCW
	Ylang EO <sup>x</sup>	Albert Vieille
	Galaxolide <sup>y</sup>	IFF
	Heliotropex <sup>z</sup>	IFF
	Phenoxyethyl isobutyrate <sup>z</sup>	Givaudan
	Isoeugenol <sup>z</sup>	Sigma Aldrich

Hedione <sup>aa</sup>	FIRMENICH
Tagetes EO <sup>aa</sup>	BLH PIM
Benzyl propionate <sup>ab</sup>	Sigma Aldrich
Sandela <sup>ab</sup>	Givaudan
Timberol <sup>ab</sup>	BLH PIM
Terpineol <sup>ac</sup>	FIRMENICH
Scentenal <sup>ac</sup>	FIRMENICH
Benzyl salicylate <sup>ac</sup>	Sigma Aldrich
Methyl salicylate <sup>ad</sup>	Robertet
Linalyl propionate <sup>ad</sup>	Givaudan
Phenoxanol <sup>ad</sup>	IFF
Grapefruit EO <sup>ae</sup>	MPE
Musk T <sup>ae</sup>	Sigma Aldrich
Menthone <sup>ae</sup>	MANE
Litsea cubeba EO <sup>af</sup>	Robertet
Lilial <sup>af</sup>	PRODASYNTH
L-Menthol <sup>af</sup>	BLH PIM
Methyl naphthyl ketone <sup>ag</sup>	BLH PIM
Ethyl vanillin <sup>ag</sup>	SOLVAY
Tonalide <sup>ag</sup>	Sigma Aldrich
Lemon petitgrain EO <sup>ah</sup>	Payan & Bertrand
Ethyl maltol <sup>ah</sup>	MPE
Ocimene <sup>ah</sup>	IFF
Delta-octalactone <sup>ai</sup>	PCW

	Gamma undecalactone <sup>ai</sup>	Sigma Aldrich
	Undecavertol <sup>ai</sup>	Givaudan
	Myrrh EO <sup>aj</sup>	Robertet
	Muscenone <sup>aj</sup>	FIRMENICH
	Aldehyde C12 MNA <sup>ak</sup>	BLH PIM
	Lemon EO <sup>ak</sup>	Albert Vieille
	Linalyl acetate <sup>ak</sup>	Sigma Aldrich
	Coumarin <sup>al</sup>	Sigma Aldrich
	Cyclogalbanate <sup>al</sup>	BLH PIM
	Virginia cedarwood EO <sup>am</sup>	Albert Vieille
	Camphor <sup>am</sup>	Sigma Aldrich
	Cyclamen aldehyde <sup>am</sup>	Givaudan
	Hexyl salicylate <sup>an</sup>	PRODASYNTH
	Black pepper EO <sup>an</sup>	Albert Vieille
	Gamma-terpinene <sup>ao</sup>	PRODASYNTH
	Gelsol <sup>ao</sup>	IFF
Categorization	Cocoa abs	Robertet
	Tonka bean abs	BLH PIM
	Broom abs	Robertet
	Hay abs	Mane
	Noble laurel EO	Payan & Bertrand
	Clary sage EO	Albert Vieille
	Parmanthema K	Firmenich
	Mimosa abs	BLH PIM
	1	1

Tobacco abs	Mane	985
Coffee abs	Robertet	986
Guaiac wood EO	PCW	987
Polysantol	Firmenich	088
Diacetyl	Sigma Aldrich	900
Glycolierral	Givaudan	989
		990

991

992

993 Table S2. Students' performance profile. List of the 14 metrics significantly modulated by 994 training that were selected from all tasks to create the performance profile of students at each 995 time measurement independently. ADJ: adjectives; BTW: between-student comparisons; REF: 996 reference or expert teacher-related metrics; sim: semantic similarity; WTH: within-student 997 comparison.

Task	Significant Metrics	
Description task	ADJ_count, NOUN_count, TTR_BTW, sim-WTH, sim-BTW, sim-REF	
Evocation task	sim-WTH, sim-BTW	
Visual recognition task	-	
Odor recognition task	Correct Rejections	
Categorization task	ARI-BTW, ARI-REF, words count, TTR-BTW, sim-BTW	

999	Table S3. RF and LSVM hyperparameters tuning. The average value (Mean), the standard
1000	deviation (SD), and the most frequently occurring value (Mode) are presented for each
1001	parameter in both models.

Model	Parameter	Mean	SD	Mode
RF	n_estimators	190	94.34	100
	max_depth	-1	0	-1
	min_samples_split	4.7	3.58	2
	min_samples_leaf	2	1.34	1
SVM	С	13.51	29.12	1
	max_iter	1000	0	1000