

An investigation of the stability of Free-Comment and Check-All-That-Apply in two consumer studies on red wines and milk chocolates

Benjamin Mahieu, Michel Visalli, Arnaud Thomas, Pascal Schlich

▶ To cite this version:

Benjamin Mahieu, Michel Visalli, Arnaud Thomas, Pascal Schlich. An investigation of the stability of Free-Comment and Check-All-That-Apply in two consumer studies on red wines and milk chocolates. Food Quality and Preference, 2021, 90, pp.104159. 10.1016/j.foodqual.2020.104159. hal-03191403

HAL Id: hal-03191403 https://hal.inrae.fr/hal-03191403

Submitted on 2 Jan 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

1 Title

- 2 An investigation of the stability of Free-Comment and Check-All-That-Apply in two
- 3 consumer studies on red wines and milk chocolates

4 Authors

- 5 Benjamin Mahieu^a, Michel Visalli^a, Arnaud Thomas^b, Pascal Schlich^a
- ⁶ ^aCentre des Sciences du Goût et de l'Alimentation, AgroSup Dijon, CNRS, INRAE,
- 7 Université Bourgogne Franche-Comté, F-21000 Dijon, France.
- ⁸ ^bSensoStat, Dijon, France.

9 Corresponding author:

- 10 Benjamin Mahieu; Pascal Schlich
- 11 <u>benjamin.mahieu@inrae.fr; pascal.schlich@inrae.fr</u>
- 12 Centre des Sciences du Goût et de l'Alimentation, AgroSup Dijon, CNRS, INRAE,
- 13 Université Bourgogne Franche-Comté, F-21000 Dijon, France.

14 Highlights

- 15 Stability of product configurations was investigated
- Stability of product by descriptor associations was investigated
- Free-Comment outputs were slightly more stable than Check-All-That-Apply
 ones
- 19 Stability was strongly dependent on the size of product differences
- Product configurations were more stable than product by descriptor
 associations

22 Abstract

- 23 Free-Comment (FC), as a response to open-ended questions, enables a word-based
- 24 sensory description and discrimination of sets of products. The stability of FC outputs
- has never been investigated and is the purpose of the present paper. Since Check-All-
- 26 That-Apply (CATA) is the most popular method for the word-based sensory description

of products with consumers, the stability of FC was compared to that of CATA 27 performed on the same products. Four red wines and four milk chocolates were 28 evaluated according to different sensory modalities by groups of consumers following 29 either an FC or a CATA protocol. The stability of the product configurations and the 30 product by descriptor associations were investigated. FC outputs were slightly more 31 stable than CATA ones. Sixty consumers enable to guarantee medium stability, if not 32 good, of FC and CATA outputs when the investigated product space is characterized 33 by large differences between the products. The minimum number of consumers to 34 obtain stable results was strongly dependent on the size of the differences between 35 the products, which suggests that if *a priori* knowledge on the size of the differences 36 37 between the investigated products is available, it must drive the decision of the number of consumers to include in the study rather than relying on an absolute rule. For both 38 39 FC and CATA, the product configurations were more easily stable in terms of numbers of consumers than the product by descriptor associations. Investigating the stability of 40 41 the product by descriptor associations a posteriori is recommended for future FC and CATA studies. 42

43 Keywords

- 44 Open-ended questions
- 45 Stability
- 46 Sensory method comparison
- 47 Consumer study

48 1. Introduction

Free-Comment (FC) (ten Kleij & Musters, 2003), as a response to open-ended 49 questions, is a sensory method that enables collecting word-based sensory 50 descriptions of a set of products without a predefined list of descriptors. For each 51 evaluated product, consumers are asked to describe the product in their own words 52 (Ares, Giménez, Barreiro, & Gámbaro, 2010; Hanaei, Cuvelier, & Sieffermann, 2015; 53 Lahne, Trubek, & Pelchat, 2014; Luc, Lê, & Philippe, 2020; Mahieu, Visalli, Thomas, 54 & Schlich, 2020; Symoneaux, Galmarini, & Mehinagic, 2012; ten Kleij & Musters, 55 2003). FC has already proven itself an efficient method in characterizing and 56 discriminating sets of products both with consumers and experts (Lahne et al., 2014; 57

Lawrence et al., 2013; ten Kleij & Musters, 2003) even out of the lab (Mahieu et al.,
2020).

Check-All-That-Apply (CATA) (Adams, Williams, Lancaster, & Foley, 2007) is a 60 61 sensory method based on a predefined list of descriptors that enables collecting wordbased sensory descriptions of sets of products. For each evaluated product, 62 consumers are asked to choose among a list of descriptors, those that apply to the 63 product. CATA also has proven itself an efficient method for the characterization and 64 discrimination of sets of products with consumers (Oppermann, de Graaf, Scholten, 65 Stieger, & Piqueras-Fiszman, 2017; Valentin, Chollet, Lelièvre, & Abdi, 2012; Varela & 66 Ares, 2012). 67

Probably because of the lack of tools for FC data analysis and ease of use of CATA, CATA is the most popular method for the word-based description of products with consumers. However, FC can provide better product discrimination as well as a richer characterization of the products as compared to CATA (Mahieu et al., 2020). Yet, while CATA has been suggested to provide stable outputs with a minimum of 60-80 consumers when differences between the products are large (Ares, Tárrega, Izquierdo, & Jaeger, 2014), the stability of the outputs provided by FC remains an open question.

In addition to the ability to characterize and discriminate the products, it is assumed that sensory methods should provide similar outputs across repeated experiments conducted in similar experimental settings. In consumer studies, it is also assumed that the larger the consumer panel, the more stable the outputs should be, but the more expensive the study is in terms of time and budget. For these reasons, having *a priori* knowledge of the number of consumers necessary to obtain stable outputs is important.

For consumer-oriented sensory methods, gathering a large number of different 82 experiments conducted under similar experimental settings with different panel sizes 83 is nearly impossible for practical limitations (Ares, Tárrega, et al., 2014). Thus, the 84 85 stability of the outputs is often evaluated internally, rather than externally, using bootstrap resampling of an actual panel that performed a study in the experimental 86 settings under interest (Ares, Bruzzone, et al., 2014; Ares, Tárrega, et al., 2014; 87 Blancher, Clavier, Egoroff, Duineveld, & Parcon, 2012; Cadena et al., 2014; 88 Mammasse & Schlich, 2014; Vidal et al., 2014; Vidal, Tárrega, Antúnez, Ares, & 89 Jaeger, 2015). This procedure enables to generate a large number of virtual panels of 90

91 different sizes that simulate repeated experiments under similar experimental settings.

92 The outputs obtained from the actual panel are considered as a benchmark to which

93 those of the virtual panels are compared.

94 Depending on the sensory method under investigation, different aspects of the outputs are compared between the actual and the virtual panels. The product configurations 95 between the actual and the virtual panels were compared in every aforementioned 96 study using the RV coefficient (Escoufier, 1973; Robert & Escoufier, 1976). For word-97 98 based sensory methods, the descriptor configurations were also compared using the RV coefficient (Ares, Bruzzone, et al., 2014; Ares, Tárrega, et al., 2014; Vidal et al., 99 100 2015). However, the descriptor configurations are usually not interpreted for themselves but rather together with the product configurations to characterize the 101 102 product space. Thus, investigating the stability of the product by descriptor associations rather than the stability of the descriptor configurations seems to be more 103 in line with common practices. 104

To the best of our knowledge, in the context of consumer word-based sensory 105 methods, no methodology has been proposed in the literature to compare the outputs 106 of the product by descriptor associations of the actual and the virtual panels. The 107 present paper proposed a methodology to do so and applied it on 10 datasets 108 corresponding to the evaluation of red wines and milk chocolates on different sensory 109 110 modalities by consumers using FC or CATA. The first objective was to investigate the number of consumers necessary to ensure the stability of FC outcomes. The second 111 objective was to compare FC and CATA conducted in similar experimental settings on 112 the stability of the outputs they provided. 113

114 2. Material and methods

115 2.1. Datasets

- 116 The information concerning the datasets used in this paper and provided across the 117 material and methods section are summarized in Table 2.
- All the data were collected using TimeSens® software (INRAE, Dijon, France).
- 119 2.1.1. First study: red wines
- 120 The datasets of this study are the same from Mahieu et al. (2020).

121 *2.1.1.1. Participants*

One-hundred and twenty consumers being 18 to 60 years old participated in this study. 122 They were recruited from a population registered in the ChemoSens Platform's 123 PanelSens database. This database has been declared to the relevant authority 124 (Commission Nationale Informatique et Libertés-CNIL-n° d'autorisation 1148039). 125 The consumers recruited were consumers of red wines at least once every two weeks 126 and were allocated in two groups of 60 consumers. The two groups were balanced in 127 terms of age repartition and gender and they were matched for consumption frequency. 128 129 The first group performed an FC task while the second group performed a CATA task. 130 Both FC and CATA were performed at home.

131 *2.1.1.2. Products*

Four commercialized French red wines from different terroirs were used. The four
 terroirs were Bordeaux, Beaujolais, Languedoc and Val de Loire.

134 2.1.1.3. FC task and datasets

For each red wine, the FC task was carried out by sensory modality in the following order: visual, olfactory, and gustatory. For each sensory modality, the following instructions were given to the consumers:

- Visual: "Describe the visual characteristics of the wine"
- Olfactory: "Describe the olfactory characteristics of the wine"
- Gustatory: "Describe the gustatory characteristics of the wine"
- 141 No particular restriction was given to the consumers on the manner of stating their 142 descriptions.
- 143 The evaluations of the red wines using FC according to the three sensory modalities
- provided three distinct datasets named FC-Wine-Vis, FC-Wine-Olf, and FC-Wine-Gus.
- 145 2.1.1.4. CATA task and datasets

For each red wine, the CATA task was carried out by sensory modality in the following
order: visual, olfactory, and gustatory. The gustatory description was presented in two

steps to the consumers: they first evaluated the basic tastes and then the aromas. For

each sensory modality, the following instruction was given to the consumers:

"Check in the subsequent list the words that apply to this wine".

The CATA lists of visual, olfactory, and gustatory descriptors were composed of 8, 10, 151 and 19 descriptors respectively. The visual descriptors were the following: violet, 152 opaque, dull, light red, bright, deep red, black, and transparent. The olfactory 153 descriptors were the following: black fruit, roasted, red fruit, green vegetable, 154 peppery/spicy, ripe fruit, animal, undergrowth, herbaceous, and woody. The gustatory 155 descriptors were the following: alcohol, slight, astringent, bitter, concentrated, 156 balanced, sweet, persistent, sour, red fruit, ripe fruit, green vegetable, black fruit, 157 roasted, peppery/spicy, herbaceous, woody, undergrowth, and animal. These 158 159 descriptors were selected according to the expertise of wine professionals, considering that they should be understandable by consumers, and were presented in a different 160 161 randomized order for each consumer but with a constant order across evaluations for 162 a given consumer.

The evaluations of the red wines using CATA according to the three sensory modalities
 provided three distinct datasets named CATA-Wine-Vis, CATA-Wine-Olf, and CATA Wine-Gus.

166 2.1.2. Second study: milk chocolates

167 *2.1.2.1. Participants*

One-hundred and forty-seven consumers being 18 to 65 years old participated in this 168 study. Seventy-seven of them were recruited from a population registered in the 169 ChemoSens Platform's PanelSens database and performed an FC task at home. The 170 remaining seventy consumers were employees of the Barry Callebaut© Company (not 171 implied in sensory and consumer research) and performed a CATA task in a dedicated 172 room at the Barry Callebaut© Company. The consumers recruited were consumers of 173 milk chocolates at least once every two weeks and were not involved in the first study. 174 The two groups were balanced in terms of age repartition and gender. 175

176 *2.1.2.2. Products*

Four milk chocolate with different recipes were used: a standard Belgian milk chocolate, a Swiss milk chocolate, a milk compound chocolate, and a protein base milk chocolate.

150

180 2.1.2.3. FC task and datasets

For each milk chocolate, the FC task was carried out by sensory modality in the following order: texture and flavor in the mouth. For each sensory modality, the following instructions were given to the consumers:

- Mouth texture: "Describe the mouth texture characteristics of the chocolate"
- Mouth flavor: "Describe the mouth flavor characteristics of the chocolate"
- 186 No particular restriction was given to the consumers on the manner of stating their 187 descriptions.
- 188 The evaluations of the milk chocolates using FC according to the two sensory 189 modalities provided two distinct datasets named FC-Choc-Tex and FC-Choc-Fla.
- 190 2.1.2.4. CATA task and datasets
- For each milk chocolate, the CATA task was carried out by sensory modality in the following order: texture and flavor in the mouth. For each sensory modality, the following instruction was given to the consumers:
- 194 "Check in the subsequent list the words that apply to this chocolate".
- The CATA lists of mouth texture and mouth flavor descriptors were composed of 8 and 6 descriptors respectively. The mouth texture descriptors were the following: hard, soft, sticky, melting, coarse, fatty, creamy texture, and mouthcoating. The mouth flavor descriptors were the following: sweet, bitter, cocoa, caramel, cereal, and milky. These descriptors were selected according to the expertise of Barry Callebaut© and were presented in a different randomized order for each consumer but with a constant order across evaluations for a given consumer.
- The evaluations of the milk chocolates using CATA according to the two sensory modalities provided two distinct datasets named CATA-Choc-Tex and CATA-Choc-Fla.
- 205 2.2. Data treatment
- 206 2.2.1. FC data treatment

- All the FC data treatments were performed using R 3.5.1 (R Core Team, 2018). The lexicon provided with IRaMuTeQ© (Ratinaud, 2014) software was used for lemmatization and part-of-speech tagging. The FC datasets were treated separately with the method described in Mahieu et al. (2020) and summarized thereafter.
- The descriptions were first cleaned, lemmatized, and filtered. Then, the words with similar meanings were grouped into latent-words relying on a chi-square-distancebased ascendant hierarchical classification.
- Among all the words and latent words, only those mentioned by at least 5% of the panel for at least one product were retained for further analysis and called descriptors thereafter. The FC lists of descriptors were composed of 8 to 20 descriptors.
- The number of times each descriptor was cited for each product was computed at the panel level. Then, the corresponding contingency table containing the citation counts of each descriptor for each product was built.
- 220 2.2.2. CATA data treatment
- The CATA datasets were treated separately and identically. The number of times each descriptor was checked for each product was computed at the panel level. Then, the corresponding contingency table containing the citation counts of each descriptor for each product was built.
- 225 2.3. Data analyses
- All analyses were performed using R 3.5.1 (R Core Team, 2018).
- 227 2.3.1. Similarity of FC and CATA outputs
- For each pair product / sensory-modality, the RV coefficient (Escoufier, 1973; Robert & Escoufier, 1976) between the configuration provided by FC and CATA was computed.
- 231 2.3.2. Size of the differences between the products
- 232 For each contingency table, the following quantity (called Cramér's Phi coefficient in
- the present paper) was computed as originally proposed by (Cramér, 1946):

234
$$\phi_C = \frac{\phi^2}{\min(r-1, c-1)}$$

with ϕ^2 the phi-square index of the contingency table, r the number of rows of the 235 contingency table, and c the number of columns of the contingency table. The phi-236 square index is equal to the sum of the eigenvalues associated with the 237 Correspondence Analysis (CA) of the contingency table. The minimum between r-1238 and c-1 is the total number of axes of this CA. Like the phi-square index itself, the 239 Cramér's Phi coefficient is a measure of the intensity of the dependence between rows 240 and columns of contingency tables. Intuitively, Cramér's Phi coefficient represents the 241 average dependence captured by one CA axis. The benefit of the Cramér's Phi 242 coefficient over the phi-square index is that it provides a measure that is comparable 243 when contingency tables are of different sizes. Cramér's Phi coefficient ranges 244 between 0 (independence) and 1 (full dependence, which corresponds to a diagonal 245 contingency table). 246

In the case of word-based sensory methods, the closer to 1 the Cramér's Phi 247 coefficient, the more dependence between products and descriptors exists in the 248 249 contingency table, and thus the more different the products are. The size of the differences between the products on a given sensory modality is estimated thanks to 250 251 the Cramér's Phi coefficient in both CATA and FC. The Cramér's Phi coefficients were compared from one dataset to another to obtain a relative ranking of the datasets in 252 253 terms of size of differences between the products. For an absolute interpretation, one can refer for example to Cohen (1988). 254

255 2.3.3. Stability of the outputs

For all computations described in this section, the configurations were obtained by CA
of the contingency tables. Principal coordinates of the products and contribution
coordinates of the descriptors were used (Castura, Antúnez, Giménez, & Ares, 2016;
Greenacre, 2013).

The stability of the descriptor configurations was not investigated (Ares, Bruzzone, et al., 2014; Ares, Tárrega, et al., 2014; Vidal et al., 2015) because they are usually not interpreted for themselves but rather as help for interpretation to understand the product configurations. In this sense, the stability of the joint product by descriptor configurations and of the product by descriptor significant associations were

investigated instead. The choice to keep two indicators (joint product by descriptor 265 configurations and product by descriptor significant associations) that seem similar is 266 deliberate. The joint product by descriptor configurations corresponds to the product 267 by descriptor insights one would draw from reading the map and/or the space resulting 268 from the CA of the contingency table. By nature, this reading is subjective and 269 approximate but has the benefit of being nuanced. The product by descriptor significant 270 associations are the black and white version of the joint product by descriptor 271 configurations and corresponds to the product by descriptor insights one would draw 272 from reading the tables as presented Mahieu et al. (2020). By their statistical-based 273 nature, the product by descriptor significant associations are objective but have the 274 275 drawback of being threshold-dependent and binary.

276 2.3.3.1. Bootstrap resampling procedure

For each dataset, different sizes of virtual panels were considered ranging from 10 to the size of the actual panel, increasing with a step of 10. For each size, 1000 virtual panels were constituted. Each virtual panel was constituted by randomly drawing subjects from the actual panel with replacement. The outputs obtained from the actual panel were considered as a benchmark to which the outputs of the virtual panels were compared.

283 2.3.3.2. Product configurations

The product configurations, i.e. the relative position of the products in relation to each other in the sensory space, were compared by computing the RV coefficient (Escoufier, 1973; Robert & Escoufier, 1976) in the full space between the product configurations of the actual and the virtual panels.

288 2.3.3.3. Joint product by descriptor configurations

To compare the joint product by descriptor configurations, i.e. the position of each product in relation to the descriptor configuration in the sensory space, the scalar products in the full space between each product vector and each descriptor vector were computed for both the actual and the virtual panels. Then, these scalar products were vectorized and the Pearson correlation coefficient was computed between the vectorized vector of scalar products of the actual panel and those of the virtual panels.

295 2.3.3.4. Product by descriptor significant associations

Fisher's exact tests per cell with a one-sided greater alternative hypothesis were conducted on each contingency table. The tests were considered significant at the α risk of 5%. These tests represent the binary statistical-based relations between each product with each descriptor.

To measure the similarity between the outputs of the tests obtained in the actual panel and each virtual panel, the Phi correlation coefficient was computed. The Phi correlation coefficient is defined as follows:

303
$$\phi = \frac{ad - bc}{\sqrt{(a+b)(c+d)(a+c)(b+d)}}$$

with "a" the number of tests that were significant in both the actual panel and the virtual panel, "b" the number of tests that were significant in the actual panel but not in the virtual panel, "c" the number of tests that were not significant in the actual panel but were in the virtual panel and "d" the number of tests that were not significant in both the actual panel and the virtual panel.

The Phi correlation coefficient is a measure of the correlation between two binary variables. It ranges between -1 and 1. A value of 0 indicates that the two variables are uncorrelated. In our case, the closer to 1 the Phi correlation coefficient, the more similar the product by descriptor significant associations were between the actual and the virtual panels.

314 *2.3.3.5. Stability of outcomes*

The reading grid was the same for all the coefficients. The stability was considered 315 good when no more than 5% of the coefficients were below 0.80. The stability was 316 considered poor when more than 5% of the coefficients were below 0.50. When the 317 stability was neither good nor poor, it was considered medium. These thresholds were 318 selected according to a common absolute value (considering that in an ideal world they 319 should be equal to one). It was necessary to achieve an objective reading of the results. 320 They were the same for the three correlation coefficients to allow for a relative 321 comparison in terms of stability of the three aspects of the outputs investigated since 322 each coefficient is comparable to the others. The proposed thresholds do not intend to 323 become "gold standards". Other thresholds might have been considered and might be 324 interesting in applications. 325

To compare the 5% quantile of the distributions of the correlation coefficients to the 326 different thresholds rather than the mean of these distributions (Ares, Bruzzone, et al., 327 2014; Ares, Tárrega, et al., 2014; Blancher et al., 2012; Cadena et al., 2014; Vidal et 328 al., 2014) is more in line to what a virtual panel drawn from the bootstrap resampling 329 of the actual panel represents. Indeed, under the hypothesis where such a virtual panel 330 represents a new study conducted in similar experimental settings, similar outputs to 331 those of the actual panel considered as a benchmark are expected from this virtual 332 panel. Thus, high correlation coefficients between the outputs of the actual and the 333 virtual panel are expected. Extended to a large number of virtual panels, this line of 334 reasoning still holds, and thus considering the entire distribution rather than its mean 335 is more in line with the bootstrap hypothesis made and with what a virtual panel 336 represents. 337

338 3. Results

339 3.1. Similarity of FC and CATA outputs

Product type	Sensory modality	RV coefficient between FC and CATA configurations		
Red wine	Visual	0.90		
Red wine	Olfactory	0.84		
Red wine	Gustatory	0.86		
Milk chocolate	Mouth texture	0.93		
Milk chocolate	Mouth flavor	0.98		

Table 1: RV coefficients between FC and CATA configurations for each pair product /
 sensory-modality

342 Overall, Table 1 shows that the RV coefficients between FC and CATA configurations

343 are high, which indicates that they provided similar product configurations.

On the detailed characterization provided by FC and CATA about the products, the reader can refer to Mahieu et al. (2020) concerning the red wines. For the milk chocolates, the characterization provided by FC and CATA were overall similar: the same sensory dimensions discriminated the products.

348 3.2. Size of the differences between the products

Dataset	Product type	Sensory modality	Sensory method	Number of products	Number of subjects	Number of descriptors	Measure of the size of the differences between the products (ϕ_c)
FC-Wine-Vis	Red wine	Visual	FC	4	60	12	0.06
CATA-Wine-Vis	Red wine	Visual	CATA	4	60	8	0.03
FC-Wine-Olf	Red wine	Olfactory	FC	4	60	14	0.05
CATA-Wine-Olf	Red wine	Olfactory	CATA	4	60	10	0.02
FC-Wine-Gus	Red wine	Gustatory	FC	4	60	20	0.07
CATA-Wine-Gus	Red wine	Gustatory	CATA	4	60	19	0.02
FC-Choc-Tex	Milk chocolate	Mouth texture	FC	4	77	10	0.17
CATA-Choc-Tex	Milk chocolate	Mouth texture	CATA	4	70	8	0.20
FC-Choc-Fla	Milk chocolate	Mouth flavor	FC	4	77	8	0.13
CATA-Choc-Fla	Milk chocolate	Mouth flavor	CATA	4	70	7	0.14

Table 2: Characteristics and measure of the size of the differences between the products for each dataset.

Table 2 summarizes the characteristics and the measures of the size of the differences 350 between the products for each dataset. For FC Cramér's Phi coefficient ranged 351 between 0.05 (Wine-Olf) and 0.17 (Choc-Tex). For CATA Cramér's Phi coefficient 352 ranged between 0.02 (Wine-Olf and Wine-Gus) and 0.20 (Choc-Tex). This suggests 353 that the size of the differences between the products differed from one product type to 354 another and from one sensory modality to another. For both FC and CATA, Cramér's 355 Phi coefficients were lower for the red wines than for the milk chocolates suggesting 356 that the size of the differences was lower between the red wines than between the milk 357 358 chocolates.

359 3.3. Stability of the outputs

360 3.3.1. Product configurations

Fig. 1 shows that good stability of the product configurations was reached for Wine-Gus, Choc-Tex, and Choc-Fla with the same minimum number of consumers with FC and CATA, respectively with 10, 10, and 20 consumers. For Wine-Vis and Wine-Olf, good stability was reached with FC with fewer consumers as compared to CATA (20 vs. 40 for Wine-Vis, 30 vs. no good stability for Wine-Olf).

Overall, the average stability of the product configurations for a given size of virtual 366 367 panels and a given pair product / sensory-modality was almost the same between FC and CATA but the minimum number of consumers required to obtain good stability of 368 the product configurations whatever the dataset was 30 for FC, and 40 for CATA 369 (except for CATA-Wine-Olf, which never reached good stability) and good stability was 370 371 reached in more datasets with FC than with CATA (5 vs. 4). For both FC and CATA, the stability of product configurations was higher for the chocolate datasets, for which 372 the size of the product differences was higher. 373

374 3.3.2. Joint product by descriptor configurations

Fig. 2 shows that whatever the method, good stability of the joint product by descriptor configurations was not reached for Wine-Olf with the actual number of consumers. For Wine-Vis and Wine-Gus, good stability was reached with FC with fewer consumers compared to CATA (40 vs. 50 for Wine-Vis, 60 vs. no good stability for Wine-Gus). For Choc-Tex and Choc-Fla, good stability was reached with FC with more consumers compared to CATA (20 vs. 10 for Choc-Tex, 30 vs. 20 for Choc-Fla).

Overall, the minimum number of consumers required to obtain good stability of the joint 381 product by descriptor configurations whatever the dataset was more than 60 382 consumers for both FC and CATA but the average stability for a given pair product / 383 sensory-modality with 60 consumers and more was slightly higher with FC than with 384 CATA for some datasets (Wine-Olf and Wine-Gus) and stability was reached in more 385 datasets with FC than with CATA (4 vs. 3). For both FC and CATA, the stability of the 386 joint product by descriptor configurations increased with the size of the product 387 differences of the datasets. For both FC and CATA, the stability of joint product by 388 descriptor configurations was higher for the chocolate datasets, for which the size of 389 the product differences was higher. 390

391 3.3.3. Product by descriptor significant associations

Fig. 3 shows that whatever the method, good stability of the product by descriptor significant associations was not reached with the actual number of consumers for all datasets and the stability was poor for the red wines datasets with the actual number of consumers. Medium stability of the product by descriptor significant associations was reached for Choc-Tex with 30 consumers for FC and 20 consumers for CATA, and for Choc-Fla with 30 consumers for FC and 50 consumers for CATA.

Overall, the minimum number of consumers required to obtain at least moderately stable product by descriptor significant associations whatever the dataset was more than 60 consumers for both FC and CATA but the average stability for a given pair product / sensory-modality was higher with FC than with CATA with 60 consumers and more for all datasets except Choc-Text. For both FC and CATA, the stability of product by descriptor significant associations was higher for the chocolate datasets, for which the size of the product differences was higher.

405 4. Discussion

406 4.1. The stability of the outputs provided by FC and CATA

Results showed relatively stable FC outputs, at least as stable as CATA ones if not
more. FC outputs reached good stability in more datasets than CATA ones regarding
product configurations and joint product by descriptor configurations. Further, the
average stability of FC outputs was always larger than or equal to CATA ones for the

three aspects of the outputs investigated in this study when a given pair product / sensory-modality with 60 consumers and more was considered. These results suggest that FC outputs are on the same level of stability that CATA ones, at least when FC and CATA are performed by sensory modality. Future studies need to be conducted to confirm or refute these results when FC and CATA are performed with a single overall characterization of each product (not by sensory modality).

The previous statements worth being nuanced by two points. First, the consumers who 417 performed the chocolate CATA task might be more knowledgeable about chocolate 418 than if they were naïve consumers. Thus, the CATA descriptions might have been 419 420 more consensual, which might have resulted in higher stability of the outputs. Therefore, the stability of CATA outputs might have been overestimated in the 421 422 chocolate study. Second, some descriptors of the CATA list in the wine study may be considered reasonably technical (e.g. animal, roasted, etc.). This may have impeded 423 the agreement of consumers on CATA descriptions, which may have resulted in lesser 424 stability of the outputs. However, some of these "technical descriptors" were mentioned 425 during the FC task (Mahieu et al., 2020), which suggests that they were meaningful to 426 consumers. They were however mentioned less frequently in FC as compared to 427 CATA, but so were common descriptors shared by FC and CATA (Mahieu et al., 2020). 428 Indeed, the CATA task encourages consumers to check the proposed descriptors 429 (Callegaro, Murakami, Tepman, & Henderson, 2015; Kim, Hopkinson, van Hout, & Lee, 430 2017; Krosnick, 1999). This suggests that this difference in citation frequency is due to 431 the task and not to the potential "technical" aspect of the descriptors. 432

Not surprisingly, for both FC and CATA, the stability of the product configurations
increased with the size of the virtual panel and with the size of the differences between
the products. The minimum number of subjects to obtain stable product configurations
was of the same order of magnitude that was previously reported for CATA, RATA,
Projective Mapping, Sorting, and Polarized Sensory Positioning (Ares, Bruzzone, et
al., 2014; Ares, Tárrega, et al., 2014; Blancher et al., 2012; Cadena et al., 2014; Vidal
et al., 2015).

The overall level of stability was more impacted by the size of product differences than by the method used (FC versus CATA). These results are in line with some previously reported studies (Ares, Bruzzone, et al., 2014; Ares, Tárrega, et al., 2014; Blancher et al., 2012; Mammasse & Schlich, 2014; Vidal et al., 2015), even with sensory descriptive analysis (Gacula Jr & Rutenbeck, 2006; Heymann, Machado, Torri, &
Robinson, 2012; Silva, Minim, Silva, & Minim, 2014). This effect of the size of product
differences affected the stability of both FC and CATA in the same direction and with
the same magnitude.

For both FC and CATA, the product configurations were more stable than the joint 448 product by descriptor configurations, themselves being more stable than the product 449 by descriptor significant associations. This suggests that the more an aspect of the 450 outputs is demanding, the less it is stable. The product configurations are relatively 451 stable because they are driven by intrinsic differences between the products and do 452 453 not depend on how these intrinsic differences are transcribed and/or verbalized. This is supported by several studies that compared two or more consumer sensory methods 454 455 and observed that they provided similar product configurations (Ares, Bruzzone, et al., 2014; Fleming, Ziegler, & Hayes, 2015; Oppermann et al., 2017; Reinbach, Giacalone, 456 Ribeiro, Bredie, & Frøst, 2014). The joint product by descriptor configurations is less 457 stable than the product configuration because identifying differences is easier than 458 explicitly verbalizing them. However, the joint product by descriptor configurations is 459 still relatively stable because the big picture of each joint product by descriptor 460 configuration is likely to be recovered across repeated experiments. The product by 461 descriptor significant associations is at best moderately stable because they require 462 the intrinsic product differences to be verbalized significantly with the same descriptors 463 across repeated experiments, which is the most demanding aspect of the outputs. 464

465 4.2. Recommendations

When the investigated product space is characterized by large differences between 466 the products, 60 consumers enable to guarantee at least a medium stability of FC and 467 CATA outputs, which is in line with previous results concerning CATA (Ares, Tárrega, 468 et al., 2014). When differences between the products are more subtle, 60 consumers 469 enable to guarantee at least a medium stability of the product configurations and the 470 joint product by descriptor configurations for both FC and CATA but do not guarantee 471 stable product by descriptor significant associations. Future studies need to be 472 conducted to investigate the number of consumers necessary to obtain good stability 473 of the product by descriptor significant associations when working with products having 474 subtle differences between them. 475

The previous recommendations are worthy of being nuanced by the fact that the 476 stability of the outputs highly depends on the size of the differences between the 477 products. Thus, these recommendations should be considered as an order of 478 magnitude rather than an absolute rule. If the practitioner has a priori knowledge of the 479 size of the differences between the products investigated, this information must be the 480 principal driver to decide the number of consumers to include in the study. Practically, 481 this a priori knowledge can arise from the relative comparison in terms of product 482 differences of the product space investigated to product spaces previously investigated 483 484 for which the stability of the outputs could have been investigated a posteriori.

Finally, like several authors recommended for the product configurations (Ares, Tárrega, et al., 2014; Blancher et al., 2012; Vidal et al., 2014), investigating *a posteriori* the stability of the joint product by descriptor configurations and of the product by descriptor significant associations is recommended to determine the degree of confidence one should have in the product by descriptor insights obtained from the study.

491 **5.** Conclusion

FC outputs were slightly more stable than CATA ones. When the product space 492 493 investigated is characterized by large differences between the products, 60 consumers enable to guarantee medium stability, if not good, of FC and CATA outputs. The 494 minimum number of consumers to obtain stable results was strongly dependent on the 495 size of the differences between the products, which suggests that if a priori knowledge 496 on the size of the differences between the products investigated is available, it must 497 drive the decision of the number of consumers to include in the study rather than an 498 absolute rule. For both FC and CATA, the sensory spaces obtained from 499 Correspondence Analysis were more stable than the product by descriptor significant 500 associations obtained from Fisher's exact tests per cell. Among sensory spaces, the 501 product configurations were more stable than the joint product by descriptor 502 503 configurations. Finally, the stability of joint product by descriptor configurations and product by descriptor significant associations are recommended to be investigated a 504 *posteriori* in the same manner that the stability of product configurations is. 505

506 Acknowledgments

- 507 This study is part of a Ph.D. financed by the Region Bourgogne-Franche-Comté and
- 508 the SensoStat Company.
- 509 The authors would like to thank Robert et Marcel®, Sicarex®, and Barry Callebaut®
- 510 for providing their products.
- 511 **References**
- Adams, J., Williams, A., Lancaster, B., & Foley, M. (2007). Advantages and uses of
 check-all-that-apply response compared to traditional scaling of attributes for
 salty snacks. In, *7th Pangborn Sensory Science Symposium*. Minneapolis,
 USA.
- Ares, G., Bruzzone, F., Vidal, L., Cadena, R. S., Giménez, A., Pineau, B., et al. (2014).
 Evaluation of a rating-based variant of check-all-that-apply questions: Rate-allthat-apply (RATA). *Food Quality and Preference, 36*, 87-95.
- Ares, G., Giménez, A., Barreiro, C., & Gámbaro, A. (2010). Use of an open-ended question to identify drivers of liking of milk desserts. Comparison with preference mapping techniques. *Food Quality and Preference, 21*(3), 286-294.
- Ares, G., Tárrega, A., Izquierdo, L., & Jaeger, S. R. (2014). Investigation of the number
 of consumers necessary to obtain stable sample and descriptor configurations
 from check-all-that-apply (CATA) questions. *Food Quality and Preference, 31*,
 135-141.
- Blancher, G., Clavier, B., Egoroff, C., Duineveld, K., & Parcon, J. (2012). A method to
 investigate the stability of a sorting map. *Food Quality and Preference, 23*(1),
 36-43.
- Cadena, R. S., Caimi, D., Jaunarena, I., Lorenzo, I., Vidal, L., Ares, G., et al. (2014).
 Comparison of rapid sensory characterization methodologies for the development of functional yogurts. *Food Research International, 64*, 446-455.
- Callegaro, M., Murakami, M. H., Tepman, Z., & Henderson, V. (2015). Yes-no answers
 versus check-all in self-administered modes. International Journal of Market
 Research, 57, 203-223.
- Castura, J. C., Antúnez, L., Giménez, A., & Ares, G. (2016). Temporal Check-All-That Apply (TCATA): A novel dynamic method for characterizing products. *Food Quality and Preference, 47*, 79-90.
- Cohen, J. (1988). Statistical Power Analysis for the Behavioral Sciences (2nd ed.).
 Hillsdale, NJ: Lawrence Erlbaum Associates, Publishers.
- 540 Cramér, H. (1946). Chapter 21. The two-dimensional case. In P. U. Press,
 541 *Mathematical Methods of Statistics*.
- 542 Escoufier, Y. (1973). Le traitement des variables vectorielles. *Biometrics, 29*, 751-760.
- Fleming, E. E., Ziegler, G. R., & Hayes, J. E. (2015). Check-all-that-apply (CATA),
 sorting, and polarized sensory positioning (PSP) with astringent stimuli. Food
 Quality and Preference, 45, 41-49.
- 546 Gacula Jr, M., & Rutenbeck, S. (2006). Sample size in consumer test and descriptive 547 analysis. *Journal of Sensory Studies, 21*(2), 129-145.
- 548 Greenacre, M. (2013). Contribution Biplots. *Journal of Computational and Graphical* 549 *Statistics, 22*(1), 107-122.
- Hanaei, F., Cuvelier, G., & Sieffermann, J. M. (2015). Consumer texture descriptions
 of a set of processed cheese. *Food Quality and Preference, 40*, 316-325.

- Heymann, H., Machado, B., Torri, L., & Robinson, A. L. (2012). How many judges
 should one use for sensory descriptive analysis? *Journal of Sensory Studies*,
 27(2), 111-122.
- Kim, I.-A., Hopkinson, A., van Hout, D., & Lee, H.-S. (2017). A novel two-step rating based 'double-faced applicability' test. Part 1: Its performance in sample
 discrimination in comparison to simple one-step applicability rating. Food
 Quality and Preference, 56, 189-200.
- 559 Krosnick, J. A. (1999). Survey research. Annual Review of Psychology, 50, 537-567.
- Lahne, J., Trubek, A. B., & Pelchat, M. L. (2014). Consumer sensory perception of cheese depends on context: A study using comment analysis and linear mixed models. *Food Quality and Preference, 32*, 184-197.
- Lawrence, G., Symoneaux, R., Maitre, I., Brossaud, F., Maestrojuan, M., & Mehinagic,
 E. (2013). Using the free comments method for sensory characterisation of
 Cabernet Franc wines: Comparison with classical profiling in a professional
 context. *Food Quality and Preference, 30*(2), 145-155.
- Luc, A., Lê, S., & Philippe, M. (2020). Nudging consumers for relevant data using Free
 JAR profiling: An application to product development. *Food Quality and Preference, 79.*
- Mahieu, B., Visalli, M., Thomas, A., & Schlich, P. (2020). Free-comment outperformed
 check-all-that-apply in the sensory characterisation of wines with consumers at
 home. Food Quality and Preference, 84.
- Mammasse, N., & Schlich, P. (2014). Adequate number of consumers in a liking test.
 Insights from resampling in seven studies. *Food Quality and Preference, 31*,
 124-128.
- Oppermann, A. K. L., de Graaf, C., Scholten, E., Stieger, M., & Piqueras-Fiszman, B.
 (2017). Comparison of Rate-All-That-Apply (RATA) and Descriptive sensory
 Analysis (DA) of model double emulsions with subtle perceptual differences.
 Food Quality and Preference, 56, 55-68.
- R Core Team. (2018). R: A language and environment for statistical computing. In.
 Vienna, Austria: R Foundation for Statistical Computing.
- 582 Ratinaud, P. (2014). IRaMuTeQ : Interface de R pour les Analyses 583 Multidimensionnelles de Textes et de Questionnaires. In. France.
- Reinbach, H. C., Giacalone, D., Ribeiro, L. M., Bredie, W. L. P., & Frøst, M. B. (2014).
 Comparison of three sensory profiling methods based on consumer perception:
 CATA, CATA with intensity and Napping[®]. *Food Quality and Preference, 32*, 160-166.
- Robert, P., & Escoufier, Y. (1976). A Unifying Tool for Linear Multivariate Statistical
 Methods: The RV- Coefficient. *Applied Statistics, 25*(3), 257-265.
- Silva, R. d. C. d. S. N. d., Minim, V. P. R., Silva, A. N. d., & Minim, L. A. (2014). Number
 of judges necessary for descriptive sensory tests. *Food Quality and Preference*,
 31, 22-27.
- Symoneaux, R., Galmarini, M. V., & Mehinagic, E. (2012). Comment analysis of
 consumer's likes and dislikes as an alternative tool to preference mapping. A
 case study on apples. *Food Quality and Preference, 24*(1), 59-66.
- ten Kleij, F., & Musters, P. A. D. (2003). Text analysis of open-ended survey responses:
 a complementary method to preference mapping. *Food Quality and Preference*,
 14(1), 43-52.
- Valentin, D., Chollet, S., Lelièvre, M., & Abdi, H. (2012). Quick and dirty but still pretty
 good: a review of new descriptive methods in food science. *International Journal* of Food Science & Technology, 47(8), 1563-1578.

- Varela, P., & Ares, G. (2012). Sensory profiling, the blurred line between sensory and
 consumer science. A review of novel methods for product characterization.
 Food Research International, 48(2), 893-908.
- Vidal, L., Cadena, R. S., Antúnez, L., Giménez, A., Varela, P., & Ares, G. (2014).
 Stability of sample configurations from projective mapping: How many consumers are necessary? *Food Quality and Preference, 34*, 79-87.
- Vidal, L., Tárrega, A., Antúnez, L., Ares, G., & Jaeger, S. R. (2015). Comparison of
 Correspondence Analysis based on Hellinger and chi-square distances to
 obtain sensory spaces from check-all-that-apply (CATA) questions. *Food Quality and Preference, 43*, 106-112.

Fig. 1: Mean of the distribution of the RV coefficients between the actual and the virtual product configurations as a function of the virtual panel size for (a) FC-Wine-Vis, (b) CATA-Wine-Vis, (c) FC-Wine-Olf, (d) CATA-Wine-Olf, (e) FC-Wine-Gus, (f) CATA-Wine-Gus, (g) FC-Choc-Tex, (h) CATA-Choc-Tex, (i) FC-Choc-Fla and (j) CATA-Choc-Fla. Dashed lines indicates 0.80 (green) and 0.50 (red). Error bars show the 0.05 and 1 quantiles of the distributions.

Fig. 2: Mean of the distribution of the Pearson correlation coefficients between the actual and the virtual joint product by descriptor configurations as a function of the virtual panel size for (a) FC-Wine-Vis, (b) CATA-Wine-Vis, (c) FC-Wine-Olf, (d) CATA-Wine-Olf, (e) FC-Wine-Gus, (f) CATA-Wine-Gus, (g) FC-Choc-Tex, (h) CATA-Choc-Tex, (i) FC-Choc-Fla and (j) CATA-Choc-Fla. Dashed lines indicates 0.80 (green) and 0.50 (red). Error bars show the 0.05 and 1 quantiles of the distributions.

Fig. 3: Mean of the distribution of the Phi correlation coefficients between the actual and the virtual Fisher's exact tests per cell ($\alpha = 5\%$) outputs as a function of the virtual panel size for (a) FC-Wine-Vis, (b) CATA-Wine-Vis, (c) FC-Wine-Olf, (d) CATA-Wine-Olf, (e) FC-Wine-Gus, (f) CATA-Wine-Gus, (g) FC-Choc-Tex, (h) CATA-Choc-Tex, (i) FC-Choc-Fla and (j) CATA-Choc-Fla. Dashed lines indicates 0.80 (green) and 0.50 (red). Error bars show the 0.05 and 1 quantiles of the distributions.





