

On the implementation of computerized adaptive observations for psychological assessment

Abstract

1
2 The use of observational tools in psychological assessment has decreased in recent years, mainly due to
3 its personnel and time costs, and researchers have not explored methodological innovations like
4 adaptive algorithms in observational assessment. In the present study, we introduce the behavior-driven
5 observation procedure to develop, test, and implement observational adaptive instruments. In Study 1,
6 we use a preexisting observational checklist to evaluate nonverbal behaviors related to psychotic
7 symptoms and to specify the adaptive algorithm's model. We fit the model to observational data
8 collected from 114 participants. The results support the model's goodness of fit. In Study 2, we use the
9 estimated model parameters to calibrate the adaptive procedure and test the algorithm for accuracy and
10 efficiency in adaptively reconstructing 58 non-adaptively collected response patterns. The results show
11 the algorithm's good accuracy and efficiency, with a 40% average reduction in the number of
12 administered items. In Study 3, we used real raters to test the adaptive checklist built with
13 behavior-driven observation. The results indicate adequate intra-rater agreement and good consistency
14 of the observed response patterns. In conclusion, the results support the possibility of using
15 behavior-driven observation to create accurate and affordable (in terms of resources) observational
16 assessment tools.

17 **Keywords:** adaptive psychological assessment; behavioral observation; behavior-driven
18 observation; one-zero sampling; modal response patterns; schizophrenia, cross-validation.

19

On the implementation of computerized adaptive observations for psychological assessment

20

21 In assessing mental disorders, the need for much high-quality information often collides with
22 the time required to collect it. The time required to entirely administer diagnostic tools, such as
23 observational checklists or semi-structured interviews (e.g., the Positive and Negative Symptoms Scale;
24 Kay, et al., 1987), sometimes could not fit the clinical routine in health-care settings (Kølbæk et al.,
25 2018), even if the tools provide extremely accurate information. On the other hand, diagnostic tools
26 such as questionnaires offer quick administration but sometimes redundant and inexhaustive
27 information, from a qualitative point of view (Serra, et al., 2015). Researchers often balance the quality
28 and quantity of information with administration time by defining short forms of assessment
29 instruments. In recent years, another increasingly frequent solution in clinical research involves
30 computerized adaptive assessment. This evaluation procedure—usually performed using electronic
31 devices such as a PC or a tablet—consists of adaptively administering items based on people’s previous
32 responses (Petersen et al., 2006; Spoto et al., 2018).

33 Clinicians have used adaptive assessments previously. In fact, this method already occurs in
34 semi-structured interviews, such as in the Structured Clinical Interview for the DSM (SCID; First,
35 2014), where the clinician moves through the questions based on the collected responses. Clinicians
36 also widely apply adaptive assessment in psychological testing via the so-called computerized adaptive
37 testing (CAT; Wainer, 2000). Several articles in the last two decades showed that CAT allows
38 clinicians to (a) reduce the length of the questionnaires (Petersen et al., 2006), (b) collect the same
39 amount of information as tests’ extended (traditional) versions (Spoto et al., 2018), (c) avoid accuracy
40 loss, and (d) increase the assessment’s efficiency (Donadello et al., 2017). While adaptive
41 questionnaires are becoming frequently adopted by researchers (Gibbons et al., 2012; Michel et al.,
42 2018; Serra, et al., 2017; Spoto et al., 2018), the same cannot be said for adaptive observational

43 assessments. Until now, the logic of adaptivity in observations is mainly related to training behavior
44 detection (e.g., the Train-to-Code software; Ray & Ray, 2008), and the detection of behaviors for
45 clinical purposes lacks sufficient research. In this sense, Pino et al. (2018) first showed how some
46 modules of the latest Autism Diagnostic Observation Schedule version (ADOS-2; C. Lord, Luyster, et
47 al., 2012; C. Lord, Rutter, et al., 2012) have good implementation potential in a computerized adaptive
48 assessment system. Unfortunately, no concrete adaptive observational tool currently exists.

49 Defining an adaptive observation involves the same critical issues encountered in constructing
50 traditional observational tools:

- 51 • the definition of a deterministic assessment model comprising a list of clear and
52 easy-to-operationalize behaviors (Hawes, et al., 2013; Haynes & O'Brien, 2000);
- 53 • the possibility of using multiple observations per person to obtain more precise response
54 patterns (observational measures require a remarkable amount of resources in terms of
55 personnel and time, so observational assessments often consist of a single scoring of an
56 instrument during [i.e., online] or after [i.e., offline] the observation, which could lack accuracy
57 because the memory interference could introduce bias); and
- 58 • the implementation and testing of the found model according to a probabilistic framework,
59 which should account for common sources of error during the observation, such as the halo
60 effect (Mumma, 2002), confirmatory bias (e.g., Cantor & Mischel, 1979), five-minute
61 impression formation (Lee, Barak, & Uhlemann, 1999), and primacy and recency effects
62 (Groth-Marnat, 2009).

63 After accomplishing these steps, observers should implement the adaptive observational tool
64 into a computerized algorithm to guide clinical observations. Finally, observers must have a report
65 containing information such as the score, the person's behavioral pattern, and the set of symptoms
66 endorsed by that pattern.

67 This methodological study aims at introducing the behavior-driven observation (BDO), a
68 method used to build computerized adaptive observational instruments that could help psychologists
69 conduct exhaustive, effective, and efficient observations. This study particularly aims at describing,
70 step-by-step, a procedure that people could use to build and test an adaptive observation. To reach this
71 goal, we applied the BDO method as a practical example to upgrade a preexisting observational
72 instrument into an adaptive tool.

73 We describe the BDO procedure in this order: after describing the methodology used to define
74 an instrument (General Methods section), we present a first study to describe how to define the
75 deterministic basis of an instrument built within the BDO and how to test it using a probabilistic model
76 of assessment (Study 1). We then use Study 1's results in Study 2 to implement an adaptive algorithm
77 of the BDO. Study 2 aims also at testing, through a simulation, the adaptive algorithm's accuracy and
78 efficiency. Finally, in a pilot study, we test the adaptive instrument with real raters in Study 3. Finally,
79 we discuss all the results, limitations, and future perspectives.

80 **General Methods**

81 In the next subsections, we describe the theoretical and methodological components of the BDO
82 and introduce a practical example to clarify the main technical issues.

83 **The Formal Psychological Assessment**

84 Researchers can use the Formal Psychological Assessment (FPA; Spoto et al., 2013)
85 methodology to build questionnaires that can exhaustively investigate the symptoms of a given
86 psychological disorder (Granziol et al., 2017). The FPA formally connects two theories of
87 mathematical psychology with clinical psychology, namely the Knowledge Space Theory (KST;
88 Doignon & Falmagne, 1999; Falmagne & Doignon, 2011) and the Formal Concept Analysis (FCA;
89 Ganter & Wille, 1999; Wille, 1982) to build assessment instruments.

90 The FPA's core aspect concerns the possibility of delineating and analyzing relationships

91 between the nonempty set A of clinical issues (i.e., diagnostic criteria or symptoms linked to a mental
 92 disorder; *attributes* in the FPA) and the nonempty set Q of items investigating those same clinical
 93 issues. The term *clinical domain* refers to the Q collection of all the items investigating a disorder.

94 When considering the evaluation of a major depressive disorder from a behavioral point of
 95 view, for example, the clinical domain comprises a set of all the possibly explored items, such as “The
 96 posture of the person points downward”. Moreover, the collection of clinical issues from common
 97 clinical practices, scientific literature, or clinical sources (e.g., the DSM-5; American Psychiatric
 98 Association [APA], 2013), referred to as the diagnosis of major depressive episode, makes up set A of
 99 attributes to investigate. The symptom “Curved posture” exemplifies a possible element in this set.
 100 Each item in the domain may investigate one or more symptoms. For instance, the previous item “The
 101 posture of the person points downward” investigates the attribute “Curved posture.”

102 The *clinical context* displays all the links between items and attributes in the form of a Boolean
 103 matrix (i.e., a binary table of 0 and 1) containing the items in rows and the attributes in columns; in a
 104 table with as many rows as the number of items and as many columns as the number of attributes, each
 105 cell contains 1 whenever item q in the row investigates the attribute a in the column (it contains 0
 106 otherwise). For instance, the item “The posture of the person points downward” (q_1) in Table 1
 107 investigates the attribute “Curved posture” (a_1); the item “Both posture and gaze of the person point
 108 downward” (q_2) investigates a_1 and the attribute “Gaze downward” (a_2). In fact, Table 1 contains a 1 in
 109 the cell corresponding to q_1 and a_1 , and both cells corresponding to a_1 and a_2 contain 1 for item q_2 .

110 [INSERT TABLE 1 HERE]

111 The clinical context allows an analysis of the *prerequisite relations* among items. An item, p ,
 112 constitutes a prerequisite for another item, q , only if the set of attributes investigated by item p
 113 comprises a subset of the attributes investigated by item q . Assuming no error, we cannot observe a
 114 positive answer to item q without observing an affirmative answer to item p . We can easily check in

115 Table 1 if item q_1 constitutes a prerequisite of item q_2 . In fact, the symptoms needed to positively
116 answer item q_2 also affirmatively answer item q_1 . Importantly, we can use the prerequisite relation in
117 adaptive assessment to infer the presence of certain symptoms without directly investigating them.

118 The clinical context and the prerequisite relation define all the admissible clinical outputs of an
119 assessment instrument built with FPA, namely the *clinical concepts*. These nonnumerical outputs
120 collect all the items endorsed by a person and the set of attributes investigated by those items. The
121 clinical context defines every clinical concept, so it follows that we can know all the clinical outputs a
122 priori from the clinical context (see Spoto et al., 2010). We call the collection of all the clinical
123 concepts the *clinical structure* (\mathcal{C}). Spoto et al. (2010, 2016) details the construction of a clinical
124 structure from a clinical context. We do not focus on these issues in the present research, so we will not
125 explain them further.

126 Figure 1 displays the clinical structure defined from Table 1's clinical context. Each concept of
127 this structure, each node on the graph, contains the set of attributes necessary and sufficient to observe
128 all the items in the concept.

129 [INSERT FIGURE 1 HERE]

130 Given the set of attributes $A = \{a_1, a_2\}$ and the set of items $Q = \{q_1, q_2\}$ (both shown in Table 1),
131 this example allows us to check all the clinical concepts defined by the clinical context. Respectively,
132 they comprise the clinical concept in which the patient does not present any symptoms (i.e., the empty
133 clinical concept, \emptyset); the clinical concept in which the patient presents only the symptom "Curved
134 posture" (i.e., the clinical concept $\{\{q_1\}; \{a_1\}\}$); and the clinical concept in which the patient presents all
135 the symptoms (i.e., the clinical concept $\{\{q_1, q_2\}; \{a_1, a_2\}\}$). Moreover, we can verify the prerequisite
136 relation among items (e.g., item q_1 constitutes a prerequisite of q_2).

137 A clinical structure like Figure 1 represents a deterministic and incomplete basis for an adaptive
138 assessment, which also requires a probabilistic framework for three main reasons. First, the observed

139 response patterns could imperfectly represent an individual’s clinical concept because of noise in the
140 assessment, represented by the *false negative* (β_q) and *false positive* (η_q) rates for each item. From an
141 observational perspective, the false negative rate shows the probability of not observing a behavior that
142 actually occurs, and the false positive rate shows the probability of reporting a behavior that has not
143 actually occurred. Second, the clinical concepts could occur with different frequencies within the
144 population; that is, each concept has its own probability, π_C , of occurrence in the population. A
145 deterministic approach implicitly assumes equiprobability for each concept. Third, a deterministic
146 clinical structure cannot predict the probability of all the clinical concepts, given people’s response
147 patterns (i.e., the observed pattern necessarily coincides with the concept). The mentioned parameters
148 (i.e., π_C for every clinical concept in the clinical structure and C , β_q and η_q for every item in the
149 clinical domain) allow researchers to define a *probabilistic clinical structure*, namely an assessment
150 model that can assign a probability of occurrence in the population to each clinical concept of the
151 structure. The probabilistic model applied to clinical structures comprises the Basic Local
152 Independence Model (BLIM; Doignon & Falmagne, 1999; Falmagne & Doignon, 1988).

153 Within the BLIM, we calculate each response pattern’s occurrence probability by multiplying
154 the conditional probability of the pattern—given that a patient exists in a given clinical concept—by
155 the probability of the clinical concept. We determine the conditional probability of the response pattern
156 given the state via the false negative and false positive rates of each item, so such error parameters
157 support the entire assessment model (formal details on the BLIM appear in Appendix S1). In fact,
158 higher error parameters yield lower reliability for the collected answers. In all the previous applications
159 of FPA (e.g., Donadello et al., 2017; Pino et al., 2018; Serra et al., 2015, 2017), a single response
160 pattern per person was sufficient to fit the BLIM and, consequently, to validate the instrument from a
161 probabilistic point of view because the instruments used comprised self-report questionnaires. As
162 mentioned before, using a single scoring in observation could reduce response patterns’ accuracy

163 because several distortions and false positives or negatives can occur. We would have to obtain a single
164 pattern derived by multiple scores. The following section describes a procedure for estimating a single
165 response pattern out of several.

166 **One-Zero Sampling and Modal Response Patterns**

167 In observational studies, people frequently use multiple observations to correctly detect or
168 quantify a behavior's frequency or duration. In particular, researchers use specific sampling strategies
169 to accurately organize these multiple data collections (Altmann, 1974; Hawes et al., 2013; Powell et al.,
170 1977). The one-zero sampling method has a long tradition in observational assessment (Goodenough,
171 1928). Originally developed for observing animal behavior, it also appears in human observation, such
172 as in the assessment of nonverbal behavior in patients diagnosed with schizophrenia (Brüne et al.,
173 2008; Troisi et al., 2007) or with depression (Geerts & Brüne, 2009).

174 The one-zero sampling method consists of two main phases: splitting the observation into n
175 equal-length samples (e.g., 15–30 s) and then checking a behavior's occurrence within each sample
176 (Martin et al., 1993). Specifically, at the end of each sample (usually announced by a beeper),
177 observers should score 1 if they observe the target behavior at any time during the sample. Thus, a
178 single score exists for each behavior via the proportion of samples in which it appears. Authors debate
179 the adequacy of such a score in the literature (e.g., Rhine & Linville, 1980; Smith, 1985). Some authors
180 discouraged its usage, arguing that it represents a biased measure of frequency or duration (Altmann,
181 1974; Dunkerton, 1981). Other authors proposed some adjustments to its computation that, while
182 maintaining the procedure's advantages, allow researchers to obtain accurate frequency and duration
183 rates based on the response patterns (Smith, 1985; Suen & Ary, 1984, 1986). In general, researchers
184 consider one-zero sampling an easy-to-apply method that allows them to observe several behaviors
185 within the same sample. Finally, one-zero scores can enhance inter- and intra-rater agreement
186 (Altmann, 1974; Rhine & Linville, 1980; Troisi et al., 2007). In this study, we propose a solution to use

187 the information of multiple one-zero samples and to organize them into a single measure: the *modal*
188 *response pattern*. We can obtain this pattern by establishing for each behavior whether occurrence or
189 non-occurrence is modal across n observational samples. Table 2 displays an example of a modal
190 response pattern:

191 [INSERT TABLE 2 HERE]

192 In this case, Items 1, 4, and 5 describe patterns observed in three out of five samples. Therefore,
193 their modal value will be 1; that is, “occurrence” takes place throughout the observation. Items 2 and 3
194 describe patterns observed during two samples; therefore, their modal value will be 0; that is,
195 “non-occurrence” takes place throughout the observation. This example has a modal response pattern
196 of $Mo = \{1,0,0,1,1\}$.

197 The modal response patterns approach conveys accurate information out of a set of response
198 patterns because the mode acts as the most appropriate central tendency index for dichotomous data,
199 such as those of one-zero sampling. Other centrality indexes could inappropriately represent
200 dichotomous data: for instance, the median and the mode in this case coincide, but the mean represents
201 the proportion of ones rather than representing a centrality indicator (Manikandan, 2011; Weisberg &
202 Weisberg, 1992). Moreover, the present study’s main objective of observation concerns establishing a
203 behavior’s occurrence or absence across n observational samples and counting its occurrences in
204 various samples does not convey this information. In fact, several authors (e.g., Dunkerton, 1981;
205 Leger, 1977; Martin et al., 1993; Powell et al., 1977) argue that using counts as a score can result in a
206 biased measure of frequency.

207 Furthermore, a modal response pattern appears closest to all other response patterns collected
208 during the observation in terms of similar responses. This property makes modal response patterns
209 good candidates for summarizing information derived from multiple observations into a single measure
210 (a detailed description of this property and the corresponding example appears in Appendix S2). After

211 reconstructing univocal information out of the patterns observed in each sample, we can implement a
212 probabilistic framework for the deterministic model of the observational assessment to define its
213 adaptive counterpart.

214 In the next sections, we will describe how to define and implement a nonadaptive deterministic
215 model of observational assessments within a probabilistic structure according to the FPA. In particular,
216 we report how we determined our sample size, all data exclusions, all manipulations, and all measures
217 to first define a nonadaptive checklist (a demonstration of the BDO method) and then test its model fit
218 and estimate its error parameters (Study 1). We then use Study 1's results to calibrate the adaptive
219 algorithm and test its accuracy and efficiency in Study 2. Finally, we test the algorithm with real raters
220 in Study 3.

221 **Study 1**

222 To develop the proposed observation method, we built a nonadaptive observational checklist via
223 the FPA. Importantly, the definition and the testing of the nonadaptive checklist should only act as a
224 demonstration of the BDO's method. In particular, experts refined preexisting lists of items and
225 attributes used to evaluate the nonverbal behavior of schizophrenia. We used both lists to define the
226 deterministic side of the assessment model. We later tested the model fit and parameter estimates of the
227 assessment model using real data. We describe all the steps below.

228 **Methods**

229 *Definition of the Nonadaptive Tool*

230 We defined the nonadaptive observational instrument starting with the list of nonverbal
231 behaviors proposed by Granzio et al. (2018) and constructed it according to the FPA. The list consisted
232 of a set of 23 items containing nonverbal behaviors related to negative symptoms frequently observed
233 in schizophrenia. We rephrased the items (e.g., shortening some items) somewhat to provide more
234 precise behavioral coding. Furthermore, we added some attributes to the original set, which

235 investigated relevant symptoms such as gaze fixation (Dowiasch et al., 2016; Gaebel, 1989). Two
236 experts in the field of schizophrenia conducted both operations starting with the original lists of items
237 and attributes.

238 The experts then defined two clinical contexts independently. The experts' agreement on this
239 task was quite high ($\kappa = .91$). After discussing the few remaining disagreements, the experts defined the
240 final clinical context, consisting of 20 items that investigated all the attributes. Tables S1 and S2 in the
241 supplemental material contain the final lists of attributes and items, respectively.

242 Starting from the context, we then constructed the corresponding clinical structure, which
243 contained 6,336 clinical concepts. To have a more convenient clinical structure size, we split the
244 clinical context into a number of subcontexts (Table 3) by clustering the items and their investigated
245 attributes around nonverbal areas of interest. In doing so, we aimed to have theoretically consistent
246 subcontexts and obtain substructures that shared the minimum possible number of investigated
247 symptoms. The most convenient solution was creating two subcontexts. One contained items that
248 investigated head and body Movements, as well as gesture and facial expressivity (Table 3a); the other
249 collected items related to prosodic features and prosocial interactive manifestations (Table 3b). The
250 former cluster, which we named *Movement*, generated a structure containing 40 concepts. The latter
251 context, which we named *ProsInt*, generated a structure containing 180 clinical concepts. These two
252 clusters shared only one attribute, "Decreased reactivity to the environment," enhancing
253 between-cluster independence.

254 [INSERT TABLE 3 HERE]

255 This division into two substructures led to a total number of clinical concepts that was easier to manage
256 in terms of model testing and parameter estimation without losing information or accuracy during
257 observational assessments.

258 *Testing the Nonadaptive Tool*

259 The testing of the deterministic structures focused on various issues: selecting samples,
260 planning a systematic observation to minimize observer bias, defining reliable modal response patterns,
261 using the modal patterns as BLIM input for testing structure fit, and estimating item error parameters.
262 Such analyses represent only one step of the BDO procedure. As specified above, a comprehensive
263 validation of the nonadaptive checklist is beyond the aim of the present study. In the next section, we
264 describe the sampling strategy for the clinical and nonclinical subsamples.

265 *Sample*

266 A total of 172 Italian participants enrolled in this study on a voluntary basis. We established the
267 sample size based on previous similar studies aimed at conducting similar estimations (Donadello et al.,
268 2017; Spoto et al., 2010). The focus of sample size definition is the number of error parameters per
269 item to estimate; in this case, we have a false positive and a false negative parameter for each item.
270 Therefore, a minimum of two participants per item is necessary. In the case at hand, we had to estimate
271 22 parameters for the ProsInt substructure (i.e., 11 items \times 2 parameters per item) and 18 parameters
272 for the Movement substructure (9 items \times 2 parameters per item). This sample size allowed us to count
273 approximately nine persons per item. We recruited a clinical subsample of 38 residential inpatients
274 with the following primary diagnoses: schizophrenia, major depressive disorder with psychotic
275 behavior and bipolar disorder with psychotic behavior. We recruited all patients diagnosed with bipolar
276 disorder during the depressive phase. All patients without a diagnosis of schizophrenia presented at
277 least one negative symptom. Patients were enrolled in three psychiatric centers in Italy: The Psychiatry
278 Unit of San Salvatore Hospital in L'Aquila; the Padova University Hospital; and the Department of
279 Clinical Neurosciences, IRCCS San Raffaele Scientific Institute in Milan.

280 Psychiatrists who were experts in the schizophrenia field ascertained diagnoses using the
281 DSM-IV-TR as nosology classification system (American Psychiatric Association [APA], 2000).
282 Patients with a diagnosis of schizophrenia received first- (~20%) or second-generation (~80%)

283 antipsychotics to treat their disorder, whereas patients with other disorders received benzodiazepines,
284 tricyclic antidepressants, and selective serotonin reuptake inhibitors (SSRIs; ~70%). Coadministration
285 of benzodiazepines was possible. Inclusion criteria for the clinical group were as follows: the presence
286 of at least one psychotic symptom; ongoing treatment with a stable dose of pharmacological therapy;
287 being a native speaker of Italian. Severe traumatic brain injury, neurological disorders, intellectual
288 disability, and alcohol or substance abuse in the past six months were the exclusion criteria. The
289 nonclinical sample comprised 134 people, mostly students, selected from the population and recruited
290 in Padova. The exclusion criteria for nonclinical sample were self-reported presence of at least one of
291 the disorders mentioned above, intellectual disability, and self-reported alcohol or substance abuse in
292 the past six months. We used a snowball sampling strategy (Goodman, 1961). The initial recruitment
293 used online advertisements. Each time participants answered the advertisement and agreed to
294 participate, we encouraged them to ask other people to participate. In this way, we increased the
295 randomness of the recruitment. Table 4 shows the demographics of the entire sample. Table S3
296 contains levels of education of all the participants, divided by group.

297 [INSERT TABLE 4 HERE]

298 Although the established sample is adequate to verify model fit (e.g., Spoto et al., 2010), the
299 proportion of patients and nonclinical individuals may be not completely adequate to represent the real
300 population. Given the prevalence of psychotic spectrum disorder (0.3–0.7% lifetime; APA, 2013), the
301 proportion of patients with a psychotic spectrum disorder diagnosis in our sample was larger than in in
302 the actual population. Nonetheless, given that the main aim of the present research was to introduce a
303 method for building, testing, and implementing an adaptive observational tool, this imbalance did not
304 represent a serious problem.

305 All participants read and filled out an informed consent form before the experimental phase.
306 The psychiatrists and psychologists explained the voluntary basis of participation and the nonintrusive

307 nature of the study very carefully. Because the study consisted of a videotaped interview, we stated
308 explicitly to the participants that they could withdraw from the interview at any time without penalty or
309 a change in the therapeutic plan. We conducted this study according to the Declaration of Helsinki and
310 obtained approval from each collaborating center's ethical committee.

311 *Stimuli and Procedure*

312 We videotaped participants during the administration of the Positive and Negative Symptoms
313 Scale (PANSS; Kay et al., 1987), a semi structured interview frequently used in psychiatric settings by
314 psychiatrists at the collaborating centers. A psychologist administered the PANSS to the nonclinical
315 group. For the patients, the psychiatrists administered the PANSS during their standard assessment
316 phases, which was less demanding on patients. In all cases, the speakers sat in front of each other with
317 a remotely controlled Sony PJ410 video camera placed on top of a 120-cm tripod and positioned
318 behind the interviewer's right shoulder. This position allowed the interviewer to record the participant's
319 whole body without the interference of the expert. The interviews lasted from 30 to 45 minutes.
320 Interviewers performed all camera operations (e.g., video extraction, charging) immediately before or
321 after each interview. This allowed the interviewers to be less distracted and reduced the participants'
322 feelings of being observed. After the interview, the experts provided the participants with a detailed
323 explanation of the interview's aim and answered any questions the participants raised.

324 Videos were then downloaded and linked to a code that was unique for each participant. We
325 extracted the video samples for the scoring phase to build the stimuli according to the one-zero
326 sampling method, and we determined the samples' number and duration empirically decided based on
327 10 pilot interviews. To provide a good tradeoff between the length of the observations and the accuracy
328 of the collected information, we selected 15 one-minute samples.

329 We created the stimuli as follows. First, we split each original video into 15 samples. A Python
330 script then randomly determined which sample to extract from the original video (Van Rossum et al.,

331 2007). The script excluded the first and last five minutes. We then edited the 15 samples, adding a
332 countdown sequence and a beeper before and after each sample, respectively. Finally, we shuffled and
333 coded this set of edited samples as a single MP4 video file¹ for use during the testing phase. We used
334 the random selection and the subsequent randomization of the 15 samples to reduce order and sequence
335 effects and minimize the risk of biases such as the first-minute impression.

336 Independent raters (one male and one female) observed and coded the 172 final videos. The
337 raters were psychiatrists selected based on their considerable expertise in mental disorders belonging to
338 the psychotic spectrum, based on the number of years they have worked with these disorders and
339 related scientific publications. In particular, both raters watched and rated the 15 samples of each video
340 on a personal computer while seated 70 cm from the screen. They received instructions not to talk or
341 have any contact with each other or with the experimenter. During the pre-experimental phase, the
342 raters received a detailed description of items along with some examples of their manifestations. We
343 also explained the scoring rule and paid particular attention to those items investigating several
344 nonverbal behaviors.

345 During the observational assessment, raters observed each video sample until the beeper's
346 sound warned them of the session's end. At that point, the video stopped, and the raters filled out the
347 nonadaptive checklist (provided in paper-and-pencil form), checking an item only if the described
348 nonverbal behaviors occurred within that sample. When the raters observed several cases per day, we
349 suggested they take a break every two observed patients and limit themselves to four observations per
350 day.

351 *Model Fitting and Parameters Estimate*

352 The one-zero sampling and modal response pattern procedures led to 344 modal response

¹ We stored all original and edited videos on hard disks located in a locked, safe place that only the project's researchers could access.

353 patterns (i.e., two patterns per patient, or, one for each rater). We obtained the final set of 172 modal
354 response patterns by combining each pair of modal patterns into a single pattern by solving the few
355 disagreements between the raters using direct discussion. Finally, we split this final set into two subsets
356 to perform a cross-validation. Specifically, we used 114 modal response patterns (i.e., 89 from the
357 nonclinical group and 25 from the clinical group) to test the structures' fit and estimate the item error
358 parameters β and η . We used the remaining 58 modal response patterns (i.e., 45 from the nonclinical
359 group and 13 from the clinical group) to perform Study 2.

360 Concerning model fitting and the parameters' estimate analyses, we used an
361 expectation-maximization algorithm (EMA; Dempster et al., 1977) implemented in MATLAB. The
362 Pearson's chi-square statistic (Falmagne & Doignon, 2011) evaluated the two structures' goodness of
363 fit to the data, and we calculated the corresponding p value using a parametric bootstrap with 5000
364 replications. We computed the p value by bootstrap due to the sparseness of the data matrices, for
365 which the asymptotic distribution of the χ^2 was not completely reliable (Reiser & Vandenberg, 1994;
366 Spoto et al., 2010). Moreover, we used the algorithm to estimate all the BLIM parameters, namely the
367 π_c for each clinical concept of the structures and the error rates (i.e., β_q and η_q) for each item on both
368 subscales. High error rates for an item indicated low reliability of the responses collected through that
369 item, meaning a possible misfit of the model.

370 In recent studies (Spoto, Stefanutti et al., 2012; 2013; Spoto, Serra, et al., 2018; Stefanutti &
371 Robusto, 2009; Stefanutti et al., 2018), error rates between 0 and .1 were low, values between .11 and
372 .2 were moderate, and values between .21 and the upper boundary of .5 were high. Appendix S3
373 contains our detailed methodological rationale for the proposed thresholding.

374 **Results**

375 In general, the 172 pairs of modal response patterns presented a very high inter-rater agreement
376 ($\kappa = .94$), and the few disagreements were easily solved. The results of the first participants' subset

377 ($n = 114$) showed a good fit of both structures to the data with adequate error parameters, as shown in
378 Table 5.

379 [INSERT TABLE 5 HERE]

380 Among the β parameters of the Movement structure, Item 5 showed a moderate value ($\beta = .13$;
381 range: $.11-.2$); Item 8 ($\beta = .36$; range: $.21-.5$) showed a high value (Table 6).

382 [INSERT TABLE 6 HERE]

383 However, the estimated η parameters were extremely small for all items, indicating the
384 probability of a false positive was particularly low. Even among items belonging to the ProsInt
385 structure, all η parameters were very low. Finally, the only item that obtained high β estimates was Item
386 16 ($\beta = .27$), which investigated a set of highly specific behaviors. In other words, with few exceptions,
387 all items presented low false positive and negative rates.

388 The results of Study 1 indicate the possibility of defining a procedure to address some critical
389 issues related to the construction of an observational instrument as first step of the BDO method. In
390 particular, a rater-friendly operationalization of the behaviors and their integration according to FPA
391 methodology resulted in the construction of a nonadaptive observational checklist. The testing
392 procedure revealed that the constructed model fit real data adequately. The missing steps are the
393 application of the model to conduct an adaptive procedure for observation and its subsequent testing in
394 the field. In Study 2, we implemented the checklist built and tested in Study 1 in an adaptive fashion by
395 defining an adaptive observation algorithm. Finally, we tested the new computerized tool's accuracy
396 and efficiency with real raters (Study 3).

397 **Study 2**

398 In the last two decades, several algorithms have implemented adaptive assessment instruments
399 starting from psychometric or mathematical frameworks. For instance, various algorithms were coded
400 according to the item response theory and mainly implemented with self-report measures (Fliege et al.,

401 2005; Gibbons et al., 2012; Michel et al., 2018). The adaptive system usually applied to FPA
402 instruments is the so-called Adaptive Testing System for Psychological Disorders (ATS-PD; Donadello
403 et al., 2017). The ATS-PD was coded within the FPA theoretical framework; therefore, it can account
404 for the deterministic and probabilistic features of a clinical structure. In other words, ATS-PD uses all
405 the parameters estimated from the application of the BLIM (i.e., probabilities of the clinical concepts
406 π_c , the false negative β_q , and the false positive η_q rates of each item) to create an adaptive instrument. In
407 this regard, the ATS-PD procedure has been tested only on self-report measures investigating
408 obsessive-compulsive disorder (Donadello et al., 2017) and major depressive episodes (Spoto et al.,
409 2018). In this study, we applied it within the BDO to implement the adaptive version of the previously
410 refined checklist.

411 **Methods**

412 *Behavior Driven Observation Algorithm*

413 An adaptive algorithm implemented using ATS-PD and here applied within the BDO is based
414 on three main rules: the questioning rule, the updating rule and the stopping rule. Appendix S4
415 describes the formal aspects of the algorithm's functioning.

416 The BDO algorithm selects items to observe according to the questioning rule. In particular, the
417 algorithm selects the most informative item. In other words, it selects the item that provides the largest
418 amount of information regardless of the collected response. Whenever two or more items are eligible or
419 observation, the algorithm selects one randomly.

420 The algorithm collects the answer to the proposed question (i.e., "yes" or "no") and then applies
421 the updating rule. This rule states that if an item receives a positive answer from the rater, the
422 probability that the final output (i.e., the clinical concept) contains that item increases, and the
423 probability decreases for those outputs that do not contain this item. This symmetric reasoning is valid
424 for a collected negative response. It is important to note that, even for the algorithm, error parameters

425 play a key role because they directly influence the extent to which an item can update the
426 aforementioned probabilities. If an item has low error parameters, it will produce a substantial
427 modification of the probability distribution of all clinical concepts (see Appendix S4). In other words,
428 the lower the error parameters for each item, the higher the reliability of the collected responses, and,
429 therefore, the higher the accuracy and efficiency of the algorithm in delineating the final output.

430 The algorithm continues asking questions and updating concepts' probabilities until it reaches a
431 stopping criterion. In particular, this stopping rule is satisfied when all items are either very unlikely to
432 be included in the final output or are almost certainly included in the final output. Therefore, none of
433 the items contributes substantially to the definition of the final output. Once the algorithm reaches this
434 criterion, it stops the assessment, stores the response pattern and its corresponding clinical concept, and
435 is ready to start a new assessment. An instrument built using the BDO method can administer multiple
436 observations. In this study, the instrument performed 15 observations in this study.² After the last
437 observational sample, the algorithm ends the entire assessment, calculates the modal response pattern,
438 and defines the corresponding modal clinical concept. In the end, the BDO's algorithm generates the
439 output, which consists of the modal response pattern M , its estimated clinical concept C_M (which
440 usually coincides with M), the list of the symptoms related to C_M , their probability estimates, and the
441 number of questions required to end the assessment. Figure 2 shows an example output provided by the
442 BDO algorithm for a single session of a simulated assessment of the Movement subscale. The first
443 section of the output contains the number of items selected by the algorithm to complete the assessment
444 out of the total number of items in the scale. In this example, the algorithm asked five of the nine items
445 in the Movement subscale (see tables S1 and S2 for details); furthermore, in this section of the output is
446 reported the number of items that received an affirmative response; in the case at hand, three items out
447 of the five asked received a positive answer. The second section of the output reports information about

² The BDO algorithm is suitable also for a single observation, if necessary.

448 the items included in the clinical concept and its associated probability value. In the case at hand, the
449 clinical concept is $C = \{2,4,5,9\}$ with a probability value of .98. Notice that, even if the algorithm
450 collected only three affirmative responses, the final clinical concept C contains four items. This is due
451 to the prerequisite relations (see General methods section) among those four items. The last section of
452 the output reports, first, the attributes that the patient may present with their associated probability;
453 second, the attributes that the patient did not prove to present, with their associated probability. In the
454 case at hand, the subject displayed three out of six attributes investigated by the Movement subscale: a
455 reduction in head movements, spontaneous movements and gesture, with high probability of
456 occurrence for each of them (i.e., 1, .99 and 1). On the other hand, the procedure suggests that the
457 attributes “reduction in facial expressivity”, “decreased reactivity to the environment”, and “rigid
458 posture” (with probability 0, .5 and 0, respectively) should not characterized the patient.

459 [INSERT FIGURE 2 HERE]

460 We coded the BDO algorithm in R (R Core Team, 2018) and implemented it in Shiny R (Chang et al.,
461 2018) for research purposes.³ In addition, we tested the algorithm’s accuracy and the efficiency, as
462 described below.

463 *Simulation Design*

464 We designed a simulation study to test the BDO algorithm’s accuracy and efficiency when
465 reproducing the original data. To accomplish these goals, we coded the algorithm to run again the
466 observational assessments used to test the nonadaptive checklist. In particular, we adaptively simulated
467 the remaining subset of 58 response patterns (i.e., obtained from their 15 observational samples) and
468 used these response patterns to define the simulated modal response patterns for each patient. Finally,
469 we compared the simulated and original modal patterns to have a measure of accuracy, and we used the

³ An R package containing the entire algorithm is in production phase. Examples of the used functions are available upon request.

470 number of suggested items as a measure of efficiency.

471 ***Outcome Measures***

472 To measure efficiency, we calculated the average number of suggested items within each
473 sample of observation and across all 15 samples. We expected an average lower number of suggested
474 items for the adaptive instrument compared to its nonadaptive counterpart's 20 items (i.e., 9 items for
475 the Movement structure and 11 for the ProsInt structure).

476 We tested accuracy by analyzing two symmetric distances. The first was the distance between
477 the modal response patterns we obtained using the nonadaptive and adaptive versions of the tool. Such
478 distances corresponded to the number of discrepant answers between the two modal response patterns
479 and represented the first index of accuracy. The second symmetric distance was the number of
480 discrepant answers between the clinical concepts (i.e., the final outputs) derived from both versions of
481 the instrument. This represented the second index of accuracy. Appendix S5 contains a formal
482 description of such indexes. We expected small distances, which indicated more similarities. Higher
483 distances corresponded to a relevant inconsistency in the generated information between the two
484 outputs (Spoto et al., 2018). We expected most of the simulated modal response patterns to be equal to
485 the original ones for both Movement and ProsInt structures.

486 **Results**

487 Table 7 displays the main results concerning the algorithm's efficiency. The algorithm
488 completed the entire assessment, suggesting checking fewer items compared to the nonadaptive version
489 of the instrument. In particular, the algorithm simulated and completed the Movement structure
490 assessment by suggesting an average of 5.1 items out of 9 ($SD = 0.3$) per observation sample. In terms
491 of overall observation, the assessment suggested only 57% of items across the 15 observational
492 samples. Concerning the ProsInt subscale, the algorithm asked an average of 7.08 items out of 11 (SD
493 $= 0.4$) per observation sample. Consequently, observation across the 15 samples was completed by

494 suggesting only 64% of items.

495 [INSERT TABLE 7 HERE]

496 Table 8 displays the distances between the modal response patterns obtained from both versions
497 of the observational tool. The BDO algorithm perfectly simulated most of the original modal response
498 patterns. In particular, for 95% of the modal response patterns simulated by the BDO, the symmetric
499 distance between the original and the simulated response patterns was zero.

500 [INSERT TABLE 8 HERE]

501 The algorithm found a limited number of modal response patterns whose symmetric distance
502 was not zero. This occurred for only three modal response patterns in both structures: for each of them,
503 the algorithm simulated a modal response pattern whose symmetric distance was equal to one item.

504 For the second accuracy index (related to clinical concepts), the symmetric distance was zero
505 for all comparisons between the concepts of the Movement structure, meaning that the modal patterns
506 calculated from both versions of the instruments converged into the same clinical concept. The same
507 result was found within the ProsInt structure, where only one comparison led to a symmetric distance
508 different from zero, namely equal to two items. In sum, the results suggested that the BDO algorithm
509 completed the assessment accurately and optimized it in real time.

510 **Study 3**

511 As a final step, we conducted a pilot study to test the tool implementation according to the BDO
512 procedure. We expected that real raters using the adaptive checklist could reach modal response
513 patterns **consistent** with those obtained through the nonadaptive checklist and save time.

514 **Methods**

515 ***Raters***

516 Four female raters, different from those of Study 1, were recruited for this last study. Their ages
517 ranged between 22 and 31 years old ($M = 26$, $SD = 3.79$). Two raters (R1 and R2) were experienced

518 psychotherapists: R1 was trained in cognitive behavioral therapy, while R2 was trained in constructivist
519 psychotherapy. The other two raters (R3 and R4) were psychology students working on their master's
520 theses.

521 *Sample*

522 A subsample of 10 patients was randomly selected from the clinical group of Study 1. All
523 patients in this sample were diagnosed with schizophrenia. This sample was equally distributed for sex,
524 and participants' ages ranged between 24 and 67 years old ($M = 45.52$, $SD = 11.87$). All patients
525 presented at least one negative symptom and were treated with first- (20%) or second-generation (80%)
526 antipsychotics. All patients gave consent for this study.

527 *Stimuli*

528 The stimuli for this study were the videos of the 10 patients included in Study 3 that were
529 collected for Study 1.

530 *Training*

531 Each rater attended a brief training session before the experimental phase. The experimenter
532 explained and discussed each item belonging to the nonadaptive checklist with the raters. Great
533 attention was devoted to items describing multiple behaviors, and items whose estimated error
534 parameters in Study 1 were moderate or high. Then, the four raters were made aware that (a) each
535 video sample was independent from the others, and (b) no time sequence was followed in the samples'
536 presentation. Finally, raters watched pilot videos (of both patients and controls) and scored the
537 nonadaptive checklist for them. Their modal response patterns were matched with those obtained from
538 Study 1 and Cohen's κ were calculated. Once an average κ of .8 was reached, the training was stopped,
539 and the experimental phase began. Four videos per rater were necessary to reach this threshold.

540 *Experimental Phase*

541 We asked each rater to watch each of the 10 videos twice, then complete the checklist once in

542 its adaptive form and once in its nonadaptive form. The order of adaptive and nonadaptive versions was
543 randomized among raters. One week was set as the inter-observational interval between the adaptive
544 and nonadaptive checklist completions. Each rater watched the video 70 cm away from a 21-inch
545 screen of an iMac 8.1 with standard resolution of 1680×1050 . During the observation, participants
546 wore headphones to avoid distraction. The experimenter instructed each rater to watch each video
547 sample until its conclusion and, only at that point, to complete the checklist and press a button to watch
548 the next sample. Raters were also asked not to interact with the experimenter. Raters were blinded to
549 the particular form of the checklist they used (i.e., adaptive or nonadaptive). Thus, to make the two
550 administrations more similar to one another, the items of the nonadaptive version of the checklist were
551 administered in a random order, which varied in every sample. To reduce the risk of biased responses
552 caused by fatigue (Haidet et al., 2009), a break was planned after the evaluation of each patient.
553 Moreover, no more than two patients per day were observed by raters.

554 **Data analysis**

555 The intraclass correlation coefficient (ICC; Fisher, 1992) was used as a measure of intra-rater
556 agreement (Koo & Li, 2016) to test whether the adaptively collected modal response patterns
557 converged with those defined via the nonadaptive checklist. The *irr* package (Gamer et al., 2012) in the
558 R statistical software was used to compute ICC and its 95% confidence interval. The time-saving
559 efficiency of the adaptive checklist was tested using a generalized linear mixed-effects model
560 (GLMM). The mean time needed to administer the checklist was used as the response variable, the
561 instrument version (adaptive vs. nonadaptive) was set as its predictor (i.e., GLMM's fixed effect), and
562 the intercept for each patient was a random factor. The *lme4* package (Bates et al., 2015) was used to
563 calculate the GLMM, while the *car* package (Fox & Weisberg, 2011) was used to obtain the *p* values.
564 The decision process leading to this analytic plan is described in Appendix S6.

565 **Results**

566 Table 9 contains the results of the efficacy and efficiency tests for each rater. In general, each
567 rater showed high intra-rater agreement, ranging from .75 to .87 ($M ICC = .82$; $SD = .05$; 95% CI =
568 [.53, .92]). Expert raters were more **consistent** and took less time to fill out the checklists.

569 [INSERT TABLE 9 HERE]

570 Regarding general efficiency, a significant effect of the instrument version emerged ($\chi^2(1) =$
571 141.18, $p < .001$). We observed that administering a nonadaptive checklist required more time on
572 average than its adaptive counterpart across the 15 samples (time difference: 6 min 28 s, $t(5) = 11.88$,
573 $p < .001$, Cohen's $d = 1.97$). This result was similar for all raters (see Table 9). These results suggest
574 that using an adaptive checklist built through a procedure such as BDO can lead to accurate data across
575 multiple observations with significant savings in time.

576 **General Discussion**

577 Observational assessment can provide experts with information that is not detectable by other
578 assessment modalities (e.g., nonverbal behavior, dynamics of interaction). In psychology and
579 psychiatry, this information is extremely useful for completing exhaustive descriptions of patients'
580 symptomatology. Nonetheless, in recent years, the rate of research and new developments in
581 observational tools for psychological assessment has decreased, probably due to difficulties related to
582 observation such as time consumption and inter-rater agreement.

583 Our methodological study attempted to introduce an adaptive observational method called
584 BDO, which succeeded in conducting a more efficient and equally accurate observational assessment,
585 compared to traditional observation methods. BDO implements the FPA methodology in the
586 observational assessment framework by using the one-zero sampling procedure for multiple
587 observations. For this paper, BDO was applied to the evaluation of nonverbal behaviors related to some
588 negative symptoms frequently observed in psychotic-spectrum disorders, mainly schizophrenia. As
589 with typical FPA applications, the BDO algorithm (i.e., the clinical concept) had appreciable clinical

590 relevance because it can be used to set a therapeutic plan targeted at a patient's specific set of
591 symptoms.

592 In its first part, the application of BDO overlaps that of FPA and consists of construction of
593 clinical context, definition of the clinical structure, and probabilistic testing based on administration of
594 an observational task in a nonadaptive fashion. In the application described in this article, two contexts
595 were built: one for prosody and prosocial behaviors, and another for Movement behaviors. These
596 contexts produced two clinical structures counting 40 Movement and 180 ProsInt clinical concepts.
597 Test results for these structures showed good fit indexes of the models to the data as well as adequate
598 values for both the β and η parameter estimates for each item. These results are encouraging because
599 they suggest that obtaining valid assessment models is possible even for observation, with low
600 probabilities of false positive or negative assessments when observing multiple behaviors. One of the
601 main advantages of defining and testing an instrument using BDO is the flexibility of this method.
602 BDO's flexibility results in a very efficient way to cope with both inter-rater disagreement on
603 item-attribute assignments and changes in the attributes selected for investigation (these could
604 correspond to either possible theoretical modifications in the definition of the disorder, or to different
605 theoretical perspectives of the experts.). Whenever disagreements among experts on item-attribute
606 assignments occur, two possibilities are provided by BDO. The first possibility refers to the solution of
607 disagreements by discussion between the experts; this is the solution adopted in our study. The second
608 possibility is the selection of the item-attribute assignment that fits the data best. On the other hand,
609 changes in the detection of the relevant attributes for a disorder can easily be implemented by slightly
610 modifying the clinical context (the Boolean matrix) and verifying the effects of the modifications on
611 model fit.

612 Furthermore, the different techniques and methods applied by the BDO algorithm yielded some
613 interesting outcomes from a clinical point of view: (a) The application of the one-zero sampling for

614 multiple observations allowed us to monitor for the occurrence of several behaviors in short intervals of
615 time, reducing the risk of primacy or recency effects; (b) the randomization of observational samples
616 could have reduced effects, such as anchoring, halo, and early impression; and (c) the collection of
617 modal response patterns from each rater gave us the chance to use a single and reliable datum for each
618 person to test the structures, overcoming the debated criticisms of one-zero patterns (Dunkerton, 1981;
619 Rhine & Linville, 1980; Smith, 1985).

620 All these advantages sum to the intrinsic advantage of the adaptive administration of the tool,
621 which solves one of the most crucial drawbacks of observation: time consumption. The last two parts
622 of BDO application consist of combining the two structures and the parameter estimates into two
623 modules (one for the ProsInt structure and one for the Movement structure) of adaptive assessment
624 tools. This is aimed at helping to accurately complete an observation and reduce its time, even with real
625 raters. In general, the time reduction is positively affected by two main factors: the size of the structure
626 and the item error parameter values. The larger the structure, the weaker the relations among the items;
627 therefore, the adaptive procedure would be less efficient. On the other hand, the lower the error
628 parameter values, the more reliable the information collected through each item administration, and,
629 therefore, the more efficient the adaptive assessment. In the applications described in this paper, the
630 parameters were low enough that in the adaptive administration, a strong reduction of the requested
631 items was observed, and a reliable reconstruction of the modal patterns was obtained. The responses'
632 convergence was corroborated even in the last pilot study, where high intra-rater agreement was
633 observed across four different raters regardless of their clinical experience. Finally, a modest reduction
634 in terms of evaluation time was observed when the adaptive checklist was used. This result is
635 reasonable considering that it is the first trial of a new procedure. Indeed, the final aim of the BDO is
636 the construction of an online adaptive observation tool. The online adaptive modality would require
637 fewer multiple observation samples, collapse the observation and evaluation phases, and take complete

638 advantage of BDO's potential. From this perspective, our results are encouraging.

639 The studies included in this paper present some limitations: For instance, the response scale of
640 the proposed checklist is binary (yes/no). The implementation of a rating scale checklist could be useful
641 for collecting information about the gravity of a specific behavior, and this is a topic of research within
642 FPA methodology. Another partial limitation is the number of observations conducted to test the
643 structures, even if the collected data are adequate for reliable estimation of the error parameters. It is
644 common practice, in any case, to assume uniform distribution of the concepts as the starting point for
645 an adaptive procedure. Because the collection of a sufficient amount of data to reliably estimate
646 concept probabilities is impossible (it would be necessary to collect more than 1,000 observations
647 conducted by at least two different raters), this solution makes application of the adaptive assessment
648 feasible, reliable, and efficient.

649 A final limitation is the number of samples per observation. While the 15 samples can solve the
650 problem related to the reliability of the response patterns, it can be argued that such a number of
651 samples is still time-demanding and should be reduced. In fact, Study 3 showed how a statistically
652 significant reduction in time occurred but did not reflect the relevant savings in terms of observable
653 items shown in the Study 2. As a first step, this work focused on the intended efficiency as a reduction
654 of the number of asked items. More effort must be made to achieve clinical time saving. Current
655 studies are focused on solving these issues and finding the minimal number of samples for obtaining
656 accurate and reliable outputs, with the goal of creating multiple adaptive online observations that could
657 almost completely solve the time consumption issue.

658 Beyond these limitations, it is important to stress some implications of our study. The BDO
659 algorithm was allowed to merge with and take advantage of several techniques and methods: the formal
660 and metric characteristics of the FPA, the use of multiple observations from the one-zero sampling
661 method, the possibility of applying all of them through modal response patterns, and the calibration of

662 adaptive systems such as the ATS-PD. As the results show, all these factors made it possible to develop
663 a procedure that could define easy-to-use adaptive observational instruments and provide clinicians
664 with help not only in conducting observations, but also with gathering sensitive information. In
665 particular, data contained within the clinical concepts are potentially the most important clinical
666 elements of an instrument defined by BDO. As shown in Figure 2, the clinical output of this checklist
667 consists of a list of several elements: the number of observed items during the entire assessment out of
668 the total number of items in the checklist, the corresponding concept within the structure and its
669 probability, and the precise list of the observed behaviors related to specific clinical symptoms, and
670 their probability of occurrence. Each of these elements can be used by clinicians during the assessment
671 phase, providing them with more exhaustive data to formulate a case. For instance, the number of items
672 observed can be used as a patient's raw score; the list of the specific symptoms can discriminate
673 between patients with the same score but different symptomatology. In this way, a specific and detailed
674 therapeutic plan can be set. Moreover, the list of symptoms can be used to focus on the core aspects to
675 be treated. Finally, the clinical output of observational instruments built with BDO can be used for two
676 other important applications. On the one hand, BDO instruments can be used to train raters because
677 their output is immediate and can be examined and discussed among trainers and trainees. On the other
678 hand, the clinical concepts can be used as an index of treatment efficacy, as the increase or decrease of
679 symptoms' number across multiple assessments can give clinicians an idea about what (and how much)
680 is changing thanks to treatment.

681 To conclude, the BDO can be used to help experts during complex observational assessments
682 by defining instruments that provide them with a clinical output that immediately conveys the specific
683 symptoms presented by patients. This information can be used to precisely formulate cases and define
684 personalized treatment strategies to help patients with specific critical conditions. Furthermore, our
685 results pave the way for other attempts at implementing observations with computerized adaptive

686 algorithms.

687

688

References

- 689 Altmann, J. (1974). Observational study of behavior: Sampling methods. *Behaviour* , 49(3-4), 227-266.
690 doi: 10.1163/156853974X00534
- 691 American Psychiatric Association [APA]. (2000). *Diagnostic and statistical manual of mental*
692 *disorders IV -Text Revision*. Washington, DC: American Psychiatric Association.
- 693 American Psychiatric Association [APA]. (2013). *Diagnostic and statistical manual of mental*
694 *disorders: DSM-5*. Washington, DC: American Psychiatric Association.
- 695 Bates, D., Maechler, M., Bolker, B. and Walker, S. (2015). Fitting Linear Mixed-Effects Models Using
696 lme4. *Journal of Statistical Software*, 67(1), 1-48. doi:10.18637/jss.v067.i01.
- 697 Brüne, M., Sonntag, C., Abdel-Hamid, M., Lehmkämpfer, C., Juckel, G., & Troisi, A. (2008).
698 Nonverbal behavior during standardized interviews in patients with schizophrenia spectrum
699 disorders. *The Journal of nervous and mental disease*, 196(4), 282-288.
- 700 Cantor, N., & Mischel, W. (1979). Prototypes in person perception. In *Advances in experimental social*
701 *psychology* (Vol. 12, pp. 3-52). Elsevier.
- 702 Chang, W., Cheng, J., Allaire, J., Xie, Y., & McPherson, J. (2018). shiny: Web application framework
703 for r [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=shiny>
704 (R package version 1.1.0)
- 705 Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via
706 the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1),
707 1-38.
- 708 Doignon, J.-P., & Falmagne, J.-C. (1999). *Knowledge spaces*. New York: Springer.
- 709 Donadello, I., Spoto, A., Sambo, F., Badaloni, S., Granzio, U., & Vidotto, G. (2017). Ats-pd: An
710 adaptive testing system for psychological disorders. *Educational and Psychological*

- 711 *Measurement*, 77(5), 792-815. doi:10.1177/0013164416652188
- 712 Dowiasch, S., Backasch, B., Einhäuser, W., Leube, D., Kircher, T., & Bremmer, F. (2016). Eye
713 movements of patients with schizophrenia in a natural environment. *European Archives of*
714 *Psychiatry and Clinical Neuroscience*, 266(1), 43-54. doi: 10.1007/s00406-014-0567-8
- 715 Dunkerton, J. (1981). Should classroom observation be quantitative? *Educational Research*, 23(2),
716 144-151.
- 717 Falmagne, J.-C., & Doignon, J.-P. (1988). A class of stochastic procedures for the assessment of
718 knowledge. *British Journal of Mathematical and Statistical Psychology*, 41(1), 1-23.
- 719 Falmagne, J.-C., & Doignon, J.-P. (2011). *Learning spaces*. New York: Springer.
- 720 First, M. B. (2014). Structured clinical interview for the dsm (scid). The encyclopedia of clinical
721 psychology, 1-6.
- 722 Fisher, R. A. (1992). Statistical methods for research workers. In *Breakthroughs in statistics* (pp.
723 66-70). Springer.
- 724 Fliege, H., Becker, J., Walter, O., Bjorner, J., Klapp, B., & Rose, M. (2005). Development of a
725 computer-adaptive test for depression (D-CAT). *Quality of Life Research*, 14(10), 2277-2291.
726 doi: 10.1007/s11136-005-6651-9
- 727 Fox, J. and Weisberg, S. (2011). *An {R} Companion to Applied Regression, Second Edition*. Thousand
728 Oaks CA: Sage. URL: <http://socserv.socsci.mcmaster.ca/jfox/Books/Companion>
- 729 Gaebel, W. (1989). Visuomotor behavior in schizophrenia. *Pharmacopsychiatry*, 22(suppl 1), 29-34.
- 730 Gamer, M., Lemon, J., & Fellows, I. (2012). *irr: Various Coefficients of Interrater Reliability*
731 *and Agreement*. R package version 0.84. <https://CRAN.R-project.org/package=irr>
- 732 Gamer, M., Lemon, J., & Fellows, I. (2012). *irr: Various Coefficients of Interrater Reliability and*
733 *Agreement*. R package version 0.84. <https://CRAN.R-project.org/package=irr>
- 734 Ganter, B., & Wille, R. (1999). *Formal concept analysis: mathematical foundations*. Berlin-Heidelberg:

- 735 Springer Verlag.
- 736 Geerts, E., & Brüne, M. (2009). Ethological approaches to psychiatric disorders: focus on depression
737 and schizophrenia. *Australian & New Zealand Journal of Psychiatry*, 43(11), 1007-1015.
- 738 Gibbons, R. D., Weiss, D. J., Pilkonis, P. A., Frank, E., Moore, T., Kim, J. B., & Kupfer, D. J. (2012).
739 Development of a computerized adaptive test for depression. *Archives of general psychiatry*,
740 69(11), 1104-1112.
- 741 Goodenough, F. L. (1928). Measuring behavior traits by means of repeated short samples. *Journal of*
742 *Juvenile Research*, 12(230), 35.
- 743 Goodman, L. A. (1961). Snowball sampling. *The annals of mathematical statistics*, 148-170.
- 744 Granzio, U., Bottesi, G., Serra, F., Spoto, A., & Vidotto, G. (2017). New perspectives on the
745 assessment of the social anxiety disorder: The formal psychological assessment. *Journal of*
746 *Evidence-Based Psychotherapies*, 17(2), 53-68.
- 747 Granzio, U., Spoto, A., & Vidotto, G. (2018). The assessment of nonverbal behavior in schizophrenia
748 through the formal psychological assessment. *International Journal of Methods in Psychiatric*
749 *Research*, 27(1). doi: 10.1002/mpr.1595
- 750 Groth-Marnat, G. (2009). *Handbook of psychological assessment*. John Wiley & Sons.
- 751 Haidet, K. K., Tate, J., Divirgilio-Thomas, D., Kolanowski, A., & Happ, M. B. (2009). Methods to
752 improve reliability of video-recorded behavioral data. *Research in nursing & health*, 32(4),
753 465-474.
- 754 Hawes, D. J., Dadds, M. R., & Pasalich, D. (2013). Observational coding strategies. *The oxford*
755 *handbook of research strategies for clinical psychology*, 120-141.
- 756 Haynes, S. N., & O'Brien, W. H. (2000). Principles and strategies of behavioral observation. In
757 *Principles and practice of behavioral assessment* (pp. 225-263). Springer.
- 758 Kay, S. R., Fiszbein, A., & Opler, L. A. (1987). The positive and negative syndromescale (panss) for

- 759 schizophrenia. *Schizophrenia bulletin*, 13(2), 261-276.
- 760 Kølbaek, P., Blicher, A. B., Buus, C. W., Feller, S. G., Holm, T., Dines, D., . . . Østergaard, S. D.
761 (2018). Inter-rater reliability of ratings on the six-item positive and negative syndrome scale
762 (panss-6) obtained using the simplified negative and positive symptoms interview (snapsi).
763 *Nordic Journal of Psychiatry*, 72(6), 431-436. (PMID: 30037286) doi:
764 10.1080/08039488.2018.1492014
- 765 Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass correlation
766 coefficients for reliability research. *Journal of chiropractic medicine*, 15(2), 155-163.
- 767 Lee, D., Barak, A., & Uhlemann, M. (1999). Forming clinical impressions during the first five minutes
768 of the counseling interview. *Psychological Reports*, 85(3 PART 1), 835-844.
- 769 Lord, C., Luyster, R., Gotham, K., & Guthrie, W. (2012). Autism diagnostic observation schedule
770 second edition (ados-2) manual (part ii): Toddler module. *Torrance, CA: Western Psychological*
771 *Services*.
- 772 Lord, C., Rutter, M., DiLavore, P., Risi, S., Gotham, K., & Bishop, S. (2012). Autism diagnostic
773 observation schedule second edition (ados-2) manual (part 1): Modules 1-4. *Torrance, CA:*
774 *Western Psychological Services*.
- 775 Martin, P., Bateson, P. P. G., & Bateson, P. (1993). Recording methods. In *Measuring behaviour: an*
776 *introductory guide* (pp. 84-100). Cambridge University Pres.
- 777 Michel, P., Baumstarck, K., Lancon, C., Ghattas, B., Loundou, A., Auquier, P., & Boyer, L. (2018).
778 Modernizing quality of life assessment: development of a multidimensional computerized
779 adaptive questionnaire for patients with schizophrenia. *Quality of Life Research*, 27(4),
780 1041-1054. doi: 10.1007/s11136-017-1553-1
- 781 Mumma, G. (2002). Effects of three types of potentially biasing information on symptom severity
782 judgments for major depressive episode. *Journal of Clinical Psychology*, 58(10), 1327-1345.

- 783 doi: 10.1002/jclp.10046
- 784 Petersen, M. A., Groenvold, M., Aaronson, N., Fayers, P., Sprangers, M., Bjorner, J. B., et al. (2006).
785 Multidimensional computerized adaptive testing of the EORTC QLQ-C30: Basic developments
786 and evaluations. *Quality of Life Research*, 15(3), 315-329.
- 787 Pino, M. C., Spoto, A., Mariano, M., Granzio, U., Peretti, S., Masedu, F., . . . Vidotto, G. (2018).
788 Formal psychological assessment for autism spectrum disorder diagnosis: A new methodology
789 to build an adaptive testing system. *The Open Psychology Journal* , 11(1), 112-122.
- 790 Powell, J., Martindale, B., Kulp, S., Martindale, A., & Bauman, R. (1977). Taking a closer look: time
791 sampling and measurement error. *Journal of Applied Behavior Analysis*, 10(2), 325-332. doi:
792 10.1901/jaba.1977.10-325
- 793 R Core Team. (2018). R: A language and environment for statistical computing [Computer software
794 manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- 795 Ray, J. M., & Ray, R. D. (2008). Train-to-code: An adaptive expert system for training systematic
796 observation and coding skills. *Behavior research methods*, 40(3), 673-693.
- 797 Reiser, M., & VandenBerg, M. (1994). Validity of the chi-square test in dichotomous variable factor
798 analysis when expected frequencies are small. *British Journal of Mathematical and Statistical*
799 *Psychology*, 47(1), 85-107. doi: 10.1111/j.2044-8317.1994.tb01026.x
- 800 Rhine, R. J., & Linville, A. K. (1980). Properties of one-zero scores in observational studies of primate
801 social behavior: The effect of assumptions on empirical analyses. *Primates*, 21(1), 111-122.
- 802 Serra, F., Spoto, A., Ghisi, M., & Vidotto, G. (2015). Formal psychological assessment in evaluating
803 depression: a new methodology to build exhaustive and irredundant adaptive questionnaires.
804 *PloS one*, 10(4), e0122131.
- 805 Serra, F., Spoto, A., Ghisi, M., & Vidotto, G. (2017). Improving major depressive episode assessment:
806 A new tool developed by formal psychological assessment. *Frontiers in psychology*, 8, 214.

- 807 Smith, P. K. (1985). The reliability and validity of one-zero sampling: misconceived criticisms and
808 unacknowledged assumptions. *British Educational Research Journal*, *11*(3), 215-220.
- 809 Spoto, A., Bottesi, G., Sanavio, E., & Vidotto, G. (2013). Theoretical foundations and clinical
810 implications of formal psychological assessment. *Psychotherapy and psychosomatics*, *82*(3),
811 197-199.
- 812 Spoto, A., Serra, F., Donadello, I., Granziol, U., & Vidotto, G. (2018). New perspectives in the
813 adaptive assessment of depression: The ATS-PD version of the QuEDS. *Frontiers in*
814 *Psychology*, *9*(JUL). doi: 10.3389/fpsyg.2018.01101
- 815 Spoto, A., Stefanutti, L., & Vidotto, G. (2010). Knowledge space theory, formal concept analysis, and
816 computerized psychological assessment. *Behavior Research Methods*, *42*(1), 342-350.
- 817 Spoto, A., Stefanutti, L., & Vidotto, G. (2012). On the unidentifiability of a certain class of skill multi
818 map based probabilistic knowledge structures. *Journal of Mathematical Psychology*, *56*(4),
819 248-255.
- 820 Spoto, A., Stefanutti, L., & Vidotto, G. (2013). Considerations about the identification of forward-and
821 backward-graded knowledge structures. *Journal of Mathematical Psychology*, *57*(5), 249-254.
- 822 Spoto, A., Stefanutti, L., & Vidotto, G. (2016). An iterative procedure for extracting skill maps from
823 data. *Behavior Research Methods*, *48*(2), 729-741. doi: 10.3758/s13428-015-0609-9
- 824 Stefanutti, L., & Robusto, E. (2009). Recovering a probabilistic knowledge structure by constraining its
825 parameter space. *Psychometrika*, *74*(1), 83-96.
- 826 Stefanutti, L., Spoto, A., & Vidotto, G. (2018). Detecting and explaining blim's unidentifiability:
827 Forward and backward parameter transformation groups. *Journal of Mathematical Psychology*,
828 *82*, 38-51.
- 829 Suen, H. K., & Ary, D. (1984). Variables influencing one-zero and instantaneous time sampling
830 outcomes. *Primates*, *25*(1), 89-94.

- 831 Suen, H. K., & Ary, D. (1986). A post hoc correction procedure for systematic errors in time-sampling
832 duration estimates. *Journal of psychopathology and behavioral assessment*, 8(1), 31-38.
- 833 Troisi, A., Pompili, E., Binello, L., & Sterpone, A. (2007). Facial expressivity during the clinical
834 interview as a predictor functional disability in schizophrenia. a pilot study. *Progress in*
835 *Neuro-Psychopharmacology and Biological Psychiatry*, 31(2),n475-481. doi:
836 10.1016/j.pnpbp.2006.11.016
- 837 Van Rossum, G., et al. (2007). Python programming language. In *Usenix annual technical conference*
838 (Vol. 41, p. 36).
- 839 Wainer, H. (2000). CATs: Whither and whence. *ETS Research Report Series*(2).
- 840 Weisberg, H., & Weisberg, H. F. (1992). *Central tendency and variability* (No. 83). Sage.
- 841 Wille, R. (1982). Restructuring lattice theory: An approach based on hierarchies of concepts. In I. Rival
842 (Ed.), *Ordered Sets* (pp. 445-470). Dordrecht: Reidel.
- 843

844

Tables

845 **Table 1**

846 *Example of Clinical Context*

q/a	a₁	a₂
q₁	1	0
q₂	1	1

847

848 *Note.* q₁ = "the person's posture points downward;" q₂ = "both posture and gaze of the person point downward;" a₁ =
 849 curved posture; a₂ = downward gaze.

850

851 **Table 2**

852 *Example Modal Response Patterns*

853

	S₁	S₂	S₃	S₄	S₅	Mo
Item 1	0	1	1	0	1	1
Item 2	1	0	1	0	0	0
Item 3	0	1	0	0	0	0
Item 4	1	0	1	1	1	1
Item 5	1	0	1	1	1	1

854 *Note.* Mo = modal response pattern. S₁-S₅ denote observational samples.

855

856 **Table 3**

857 *Clinical Subcontexts for Movement and ProsInt Scales*

858

Movement							ProsInt									
Item	A ₁	A ₂	A ₇	A ₈	A ₁₀	A ₁₂	Item	A ₃	A ₄	A ₅	A ₆	A ₉	A ₁₀	A ₁₁	A ₁₃	A ₁₄
1	0	0	0	0	0	1	11	1	0	0	0	0	0	0	0	0
2	0	0	1	0	0	0	12	0	1	0	0	0	0	0	0	0
3	0	0	1	0	1	1	13	0	1	1	0	0	0	0	0	0
4	0	0	0	1	0	0	14	0	0	0	1	0	0	0	0	0
5	0	0	1	1	0	0	15	0	1	1	1	0	0	0	0	0
6	1	0	0	0	0	0	16	0	1	0	1	0	0	0	0	0
7	1	0	0	1	0	0	17	0	0	0	0	1	1	0	0	0
8	0	1	1	1	0	0	18	1	0	0	0	1	1	0	0	0
9	1	0	0	0	1	0	19	0	0	0	0	0	0	1	0	0
10	0	1	0	0	0	0	20	0	0	0	0	0	1	0	1	0
3a							21	0	0	0	0	0	1	0	0	1
							3b									

859

860 **Table 4**861 *Demographic Characteristics*

Group	<i>n</i>	Age Range		Sex	Occupation	Marital Status	
		<i>M ± 1 SD</i>					
Nonclinical	134	26	3.40	19-67	F: 100 M: 34	Student: 125 Working student: 7 Employed: 1 Retired: 1	Single: 130 Married: 3 Divorced: 1
Clinical	38	42	7.35	24-67	F: 10 M: 28	Student: 1 Unemployed: 30 Employed: 6 Retired: 1	Single: 24 Married: 14
Diagnosis: SCH	25	44.5	13	24-67	F: 5 M: 20	Student: 1 Unemployed: 22 Employed: 1 Retired: 1	Single: 18 Married: 7
Diagnosis: MDD	5	51	6.5	46-63	F: 3 M: 2	Unemployed: 2 Employed: 3	Single: 1 Married: 4
Diagnosis: BD	8	36	6.5	27-45	F: 2 M: 6	Unemployed: 6 Employed: 2	Single: 5 Married: 3

862

863 *Note.* SCH = schizophrenia; MDD = major depressive disorder; BD = bipolar disorder. Descriptors of the clinical group are
864 further split across diagnoses.

865

866

867 **Table 5**

868 *Model Fit Indexes for Movement and ProsInt Structures*

869

870	Structure	χ^2	<i>df</i>	Bootstrap <i>p</i>
871	Movement	43.44	454	.06
872	ProsInt	23.10	1846	.11

873

874 **Table 6**

875 *Error Parameter Estimates for Movement and ProsInt Structure Items*

876

Movement			ProsInt		
Item	β	η	Item	β	η
1	.1	.01	10	<.001	<.001
2	<.001	<.001	11	<.001	<.001
3	.1	.01	12	<.001	<.001
4	<.001	<.001	13	<.001	.052
5	.133	<.001	14	<.001	<.001
6	<.001	.068	15	.273	<.001
7	.076	<.001	16	<.001	<.001
8	.364	<.001	17	<.001	<.001
9	<.001	<.001	18	.1	.01
			19	.1	.01
			20	.1	.01

877 *Note.* β = false negative; η = false positive.

878

879 **Table 7**

880 *Number of Suggested Items in Nonadaptive and Adaptive Versions of Observational Tool*

881

Efficiency			
Structure	Instrument	<i>n</i>	<i>N</i>
Movement	Nonadaptive	9	135
Movement	Adaptive	5.10	76.50
ProsInt	Nonadaptive	11	165
ProsInt	Adaptive	7.08	106.20

882 *Note.* *n* = mean number of suggested items per sample; *N* = mean number of suggested items across 15 observation
 883 samples.

884

885

886 **Table 8**

887 *Accuracy of BDO Algorithm*

888

Accuracy			
Structure	Δ	d = 0*	d = 1**
Movement	0.05	55	3
ProsInt	0.05	55	3

889 *Note.* Δ = average symmetric distance between original and simulated modal response patterns.

890 *Number of modal simulated patterns at distance 0 from original.

891 **Number of simulated modal patterns at distance 1 from original.

892

893 **Table 9**894 *Accuracy and Efficiency of BDO Adaptive Tool, Real Raters*

895

Rater	Mean ICC	95% CI	Nonadaptive Minutes*	Adaptive Minutes	Difference
R1	.87	.60, .92	17.30 (2.37)	12.55 (3.31)	4.55 (3.28)
R2	.83	.60, .92	17.09 (3.04)	10.46 (3.04)	6.55 (2.32)
R3	.85	.49, .89	32.36 (4.18)	24.05 (3.16)	8.07 (3.01)
R4	.75	.44, .88	32.28 (5.16)	24.20 (2.36)	7.43 (6.33)

896 *Note.* *Ms* and *CI*s of ICCs across 10 patients for each rater.897 *In the last three columns, the average (\pm *SD*) amount of time required for the completion of the nonadaptive and adaptive
898 forms are reported, together with their average (\pm *SD*) differences among the 10 patients.

899

900