

1 **Modelling personality, plasticity and predictability in**
2 **shelter dogs**

3

4 Conor Goold^{1*} and Ruth C. Newberry¹

5 ¹Department of Animal and Aquacultural Sciences, Faculty of Biosciences,
6 Norwegian University of Life Sciences

7

8 *Corresponding author: conor.goold@nmbu.no

9 **Abstract**

10 Behavioural assessments of shelter dogs (*Canis lupus familiaris*) typically comprise
11 standardised test batteries conducted at one time point but test batteries have shown
12 inconsistent predictive validity. Longitudinal behavioural assessments offer an
13 alternative. We modelled longitudinal observational data on shelter dog behaviour using
14 the framework of behavioural reaction norms, partitioning variance into personality (i.e.
15 inter-individual differences in behaviour), plasticity (i.e. individual differences in
16 behavioural change) and predictability (i.e. individual differences in residual intra-
17 individual variation). We analysed data on 3,263 dogs' interactions (N = 19,281) with
18 unfamiliar people during their first month after arrival at the shelter. Accounting for
19 personality, plasticity (linear and quadratic trends) and predictability improved the
20 predictive accuracy of the analyses compared to models quantifying personality and/or
21 plasticity only. While dogs were, on average, highly sociable with unfamiliar people and
22 sociability increased over days since arrival, group averages were unrepresentative of all
23 dogs and predictions made at the individual level entailed considerable uncertainty.
24 Effects of demographic variables (e.g. age) on personality, plasticity and predictability
25 were observed. Behavioural repeatability was higher one week after arrival compared to
26 arrival day. Our results highlight the value of longitudinal assessments on shelter dogs
27 and identify measures that could improve the predictive validity of behavioural
28 assessments in shelters.

29

30 **Keywords**— inter- and intra-individual differences, behavioural reaction norms,
31 behavioural repeatability, longitudinal behavioural assessment, human-animal
32 interactions.

33

34 **Introduction**

35 *Personality*, defined by inter-individual differences in average behaviour, represents just
36 one component of behavioural variation of interest in animal behaviour research.

37 Personality frequently describes less than 50% of behavioural variation in animal
38 personality studies [1,2], leading to the combined analysis of personality with *plasticity*,
39 individual differences in behavioural change [3], and *predictability*, individual
40 differences in residual intra-individual variability [4–8]. Understanding these different
41 sources of behavioural variation simultaneously can be achieved using the general
42 framework of behavioural reaction norms [3,5], which provides insight into how animals
43 react to fluctuating environments through time and across contexts. The concept of
44 behavioural reactions norms is built upon the use of hierarchical statistical models to
45 quantify between- and within-individual variation in behaviour, following methods in
46 quantitative genetics [3]. More generally, these developments reflect increasing interest
47 across biology in expanding the ‘trait space’ of phenotypic evolution [9] beyond mean
48 trait differences and systematic plasticity across environmental gradients to include
49 residual trait variation (e.g. developmental instability: [10,11]; stochastic variation in
50 gene expression: [12]).

51

52 Modest repeatability of behaviour has been documented in domestic dogs (*Canis lupus*
53 *familiaris*), providing evidence for personality variation. For instance, using meta-
54 analysis, Fratkin *et al.* [13] found an average Pearson’s correlation of behaviour through
55 time of 0.43, explaining 19% of the behavioural variance between successive time points
56 (where the average time interval between measurements was 21 weeks). However, the
57 goal of personality assessments in dogs is often to predict an individual dog’s future
58 behaviour (e.g. working dogs: [14,15]; pet dogs: [16]) and, thus, it is important not to
59 confuse the stability of an individual’s behaviour relative to the behaviour of others with
60 stability of intra-individual behaviour. That is, individuals could vary their behaviour in
61 meaningful ways in response to internal (e.g. ontogeny) and external (e.g. environmental)
62 factors while maintaining differences from other individuals. When time-related change
63 in dog behaviour has been taken into account, behavioural change at the group-level has

64 been of primary focus (e.g. [16–18]) and no studies have explored the heterogeneity of
65 residual variance within each dog. The predominant focus on inter-individual differences
66 and group-level patterns of behavioural change risks obscuring important individual-level
67 heterogeneity and may partly explain why a number of dog personality assessment tools
68 have been unreliable in predicting future behaviour [14–16,19].

69

70 Of particular concern is the low predictive value of shelter dog assessments for predicting
71 behaviour post-adoption [20–24], resulting in calls for longitudinal, observational models
72 of assessment [20,24]. Animal shelters are dynamic environments and, for most dogs,
73 instigate an immediate threat to homeostasis as evidenced by heightened hypothalamic-
74 pituitary-adrenal axis activity and an increase in stress-related behaviours (e.g. [25–28]).
75 Over time, physiological and behavioural responses are amenable to change [17,27,29].
76 Therefore, dogs in shelters may exhibit substantial heterogeneity in intra-individual
77 behaviour captured neither by standardised behavioural assessments conducted at one
78 time point [24] nor by group-level patterns of behavioural change. An additional
79 complication is that the behaviour in shelters may not be representative of behaviour
80 outside of shelters. For example, Patronek and Bradley [29] suggested that up to 50% of
81 instances of aggression expressed while at a shelter are likely to be false positives. Such
82 false positives may be captured in estimates of predictability, with individuals departing
83 more from their representative behaviour having higher residual intra-individual
84 variability (lower predictability) than others. Overall, absolute values of behaviour, such
85 as mean trait values across time (i.e. personality), may account for just part of the
86 important behavioural variation needed to understand and predict shelter dog behaviour.
87 While observational models of assessment have been encouraged, methods to
88 systematically analyse longitudinal data collected at shelters into meaningful formats are
89 lacking.

90

91 In this paper, we demonstrate how the framework of behavioural reaction norms can be
92 used to quantify inter- and intra-individual differences in shelter dog behaviour. To do so,

93 we employ data on dogs' interactions with unfamiliar people from a longitudinal and
94 observational shelter assessment. As a core feature of personality assessments, how
95 shelter dogs interact with unknown people is of great importance. At one extreme, if dogs
96 bite or attempt to bite unfamiliar people, they are at risk of euthanasia [29]. At the other
97 extreme, even subtle differences in how dogs interact with potential adopters can
98 influence adoption success [30]. Importantly, neither may all dogs react to unfamiliar
99 people in the same way through time at the shelter nor may all dogs show the same day-
100 to-day fluctuation of behaviour around their average behavioural trajectories. These
101 considerations can be explored by examining behavioural reaction norms.

102

103 The analysis of behavioural reaction norms is dependent on the use of hierarchical
104 statistical models for partitioning variance among individuals [3,5,6]. Given that ordinal
105 data are common in behavioural research, here, we illustrate how similar hierarchical
106 models can be applied to ordinal data using a Bayesian framework (see also [31]). Apart
107 from distinguishing inter- from intra-individual variation, we place particular emphasis
108 on two desirable properties of the hierarchical modelling approach taken here. First, the
109 property of *hierarchical shrinkage* [32] offers an efficacious way of making inferences
110 about individual-level behaviour when data are highly unbalanced and potentially
111 unrepresentative of a dog's typical behaviour. When data are sparse for certain
112 individuals, hierarchical shrinkage means that an individual's parameter estimates (e.g.
113 intercepts) are more similar to, or shrunken, towards the group-level estimates. Secondly,
114 since any prediction of future (dog) behaviour will entail uncertainty, a Bayesian
115 approach is attractive because we can directly obtain a probability distribution of
116 parameter values consistent with the data (i.e. the posterior distribution) for all
117 parameters [32,33]. By contrast, frequentist confidence intervals are not posterior
118 probability distributions and, thus, their interpretation is more challenging when a goal is
119 to understand uncertainty in parameter estimates [32].

120

121 **Material & Methods**

122 **Subjects**

123 Behavioural data on $N = 3,263$ dogs from Battersea Dogs and Cats Home's longitudinal,
124 observational assessment model were used for analysis. The data concerned all
125 behavioural records of dogs at the shelter during 2014 (including those arriving in 2013
126 or departing in 2015), filtered to include all dogs: 1) at least 4 months of age (to ensure
127 all dogs were treated similarly under shelter protocols, e.g. vaccinated so eligible for
128 walks outside and kennelled in similar areas), 2) with at least one observation during the
129 first 31 days since arrival at the shelter, and 3) with complete data for demographic
130 variables to be included in the formal analysis (Table 1). Because dogs spent
131 approximately one month at the shelter on average (Table 1), we focused on this period in
132 our analyses (arrival day 0 to day 30). We did not include breed characterisation due to
133 the unreliability of using appearance to attribute breed type to shelter dogs of uncertain
134 heritage [34].

135

136 **Shelter environment**

137 Details of the shelter environment have been presented elsewhere [35]. Briefly, the
138 shelter was composed of three different rehoming centres (Table 1): one large inner-city
139 centre based in London (approximate capacity: 150-200 dogs), a medium-sized
140 suburban/rural centre based in Old Windsor (approximate capacity: 100-150 dogs), and a
141 smaller rural centre in Brands Hatch (approximate capacity: 50 dogs). Dogs considered
142 suitable for adoption were housed in indoor kennels (typically about 4m x 2m, with a
143 shelf and bedding alcove; see also [36]). Most dogs were housed individually, and given
144 daily access to an indoor run behind their kennel. Feeding, exercising and kennel
145 cleaning were performed by a relatively stable group of staff members. Dogs received

146 water ad libitum and two meals daily according to veterinary recommendations. Sensory
147 variety was introduced daily (e.g. toys, essential oils, classical music, access to quiet
148 ‘chill-out’ rooms). Regular work hours were from 0800 h to 1700 h each day, with public
149 visitation from 1000 h to 1600 h. Dogs were socialised with staff and/or volunteers daily.

150

151 **Data collection**

152 The observational assessment implemented at the shelter included observations of dogs
153 by trained shelter employees in different, everyday contexts, each with its own qualitative
154 ethogram of possible behaviours. Shortly after dogs were observed in relevant contexts,
155 employees entered observations into a custom, online platform using computers located
156 in different housing areas. Each behaviour within a context had its own code. Previously,
157 we have reported on aggressive behaviour across contexts [35]. Here, we focus on
158 variation in behaviour in one of the most important contexts, ‘Interactions with
159 unfamiliar people’, which pertained to how dogs reacted when people with whom they
160 had never interacted before approached, made eye contact, spoke to and/or attempted to
161 make physical contact with them. For the most part, this context occurred outside of the
162 kennel, but it could also occur if an unfamiliar person entered the kennel. Observations
163 could be recorded by an employee meeting an unfamiliar dog, or by an employee
164 observing a dog meeting an unfamiliar person. Different employees could input records
165 for the same dog, and employees could discuss the best code to describe a certain
166 observation if required.

167

168 Behavioural observations in the ‘Interactions with unfamiliar people’ context were
169 recorded using a 13-code ethogram (Table 2). Each behavioural code was subjectively
170 labelled and generally defined, providing a balance between behavioural rating and
171 behavioural coding methodologies. The ethogram represented a scale of behavioural
172 problem severity and assumed adoptability (higher codes indicating higher severity of
173 problematic behaviour/lower sociability), reflected by grouping the 13 codes further into
174 green, amber and red codes (Table 2). Green behaviours posed no problems for adoption,

175 amber behaviours suggested dogs may require some training to facilitate successful
176 adoption but did not pose a danger to people or other dogs, and red behaviours suggested
177 dogs needed training or behavioural modification to facilitate successful adoption and
178 could pose a risk to people or other dogs. A dog's suitability for adoption was, however,
179 based on multiple behavioural observations over a number of days. When registering an
180 observation, the employee selected the highest code in the ethogram that was observed on
181 that occasion (i.e. the most severe level of problematic behaviour was given priority).
182 There were periods when a dog could receive no entries for the context for several days
183 but other times when multiple observations were recorded on the same day, usually when
184 a previous observation was followed by a more serious behavioural event. In these
185 instances, and in keeping with the shelter protocol, we retained the highest (i.e. most
186 severe) behavioural code registered for the context that day. When the behaviours were
187 the same, only one record was retained for that day. This resulted in an average of 5.9
188 (SD = 3.7; range = 1 to 22) records per dog on responses during interactions with
189 unfamiliar people while at the shelter. For dogs with more than one record, the average
190 number of days between records was 2.8 (SD = 2.2; range = 1 to 29).

191

192 **Validity & inter-rater reliability**

193 Inter-rater reliability and the validity of the assessment methodology were evaluated
194 using data from a larger research project at the shelter. Videos depicting different
195 behaviours in different contexts were filmed by canine behaviourists working at the
196 shelter, who subsequently organised video coding sessions with 93 staff members (each
197 session with about 5 - 10 participants) across rehoming centres [35]. The authors were
198 blind to the videos and administration of video coding sessions. The staff members were
199 shown 14 videos (each about 30 s long) depicting randomly-selected behaviours, two
200 from each of seven different assessment contexts (presented in a pseudo-random order,
201 the same for all participants). Directly after watching each video, they individually
202 recorded (on a paper response form) which ethogram code best described the behaviour
203 observed in each context. Two videos depicted behaviour during interactions with people

204 (familiar versus unfamiliar not differentiated), one demonstrating *Reacts to people*
205 *aggressive* and the other *Reacts to people non-aggressive* (Table 2). Below, we present
206 the inter-rater reliabilities and the percentage of people who chose the correct behaviour
207 and colour category for these two videos in particular, but also the averaged results across
208 the 14 videos, since there was some redundancy between ethogram scales across
209 contexts.

210

211 **Statistical analyses**

212 All data analysis was conducted in R version 3.3.2 [37].

213

214 **Validity & inter-rater reliability**

215 Validity was assessed by calculating the percentage of people answering with the correct
216 ethogram code/code colour for each video. Inter-rater reliability was calculated for each
217 video using the consensus statistic [38] in the R package *agrmt* [39], which is based on
218 Shannon entropy and assesses the amount of agreement in ordered categorical responses.
219 A value of 0 implies complete disagreement (i.e. responses equally split between the
220 lowest and highest ordinal categories, respectively) and a value of 1 indicates complete
221 agreement (i.e. all responses in a single category). For the consensus statistic, 95%
222 confidence intervals (CIs) were obtained using 10,000 non-parametric bootstrap samples.
223 The confidence intervals were subsequently compared to 95% CIs of 10,000 bootstrap
224 sample statistics from a null uniform distribution, which was created by: 1) selecting the
225 range of unique answers given for a particular video and 2) taking 10,000 samples of the
226 same size as the real data, where each answer had equal probability of being chosen.
227 Thus, the null distribution represented a population with a realistic range of answers, but
228 had no clear consensus about which category best described the behaviour. When the null
229 and real consensus statistics' 95% CIs did not overlap, we inferred statistically significant
230 consensus among participants.

231

232 **Hierarchical Bayesian ordinal probit model**

233 The distribution of ethogram categories was heavily skewed in favour of the green codes
234 (Table 2), particularly the first *Friendly* category. Since some categories were chosen
235 particularly infrequently, we aggregated the raw responses into a 6-category scale: 1)
236 *Friendly*, 2) *Excitable*, 3) *Independent*, 4) *Submissive*, 5) *Amber codes*, 6) *Red codes*.
237 This aggregated scale retained the main variation in the data and simplified the data
238 interpretation. We analysed the data using a Bayesian ordinal probit model (described in
239 [32,40]), but extended to integrate the hierarchical structure of the data, including
240 heteroscedastic residual standard deviations, to quantify predictability for each dog (for
241 related models, see [31,41,42]). The ordinal probit model, also known as the cumulative
242 or thresholded normal model, is motivated by a latent variable interpretation of the
243 ordinal scale. That is, an ordinal dependent variable, Y , with categories K_j , from $j = 1$ to
244 J , is a realisation of an underlying continuous variable divided into thresholds, θ_c , for
245 $c = 1$ to $J - 1$. Under the probit model, the probability of each ordinal category is equal
246 to its area under the cumulative normal distribution, Φ , with mean, μ , SD σ and
247 thresholds θ_c :

$$Prob(Y = K | \mu, \sigma, \theta_c) = \Phi\left[\frac{\theta_c - \mu}{\sigma}\right] - \Phi\left[\frac{\theta_{c-1} - \mu}{\sigma}\right] \quad (1)$$

248 For the first and last categories, this simplifies to $\Phi[(\theta_c - \mu)/\sigma]$ and $1 - \Phi[(\theta_{c-1} -$
249 $\mu)/\sigma]$, respectively. As such, the latent scale extends from $\pm\infty$. Here, the ordinal
250 dependent variable was a realisation of the hypothesised continuum of ‘insociability
251 when meeting unfamiliar people’, with 6 categories and 5 threshold parameters. While
252 ordinal regression models usually fix the mean and SD of the latent scale to 0 and 1 and
253 estimate the threshold parameters, we fixed the first and last thresholds to 1.5 and 5.5
254 respectively, allowing for the remaining thresholds, and the mean and SD, to be estimated
255 from the data. As explained by Kruschke [32], this allows for the results to be
256 interpretable with respect to the ordinal scale. We present the results using both the

257 predicted probabilities of ordinal sociability codes and estimates on the latent,
258 unobserved scale assumed to generate the ordinal responses.

259

260 **Hierarchical structure**

261 To model inter- and intra-individual variation, a hierarchical structure for both the mean
262 and SD was specified. That is, parameters were included for both group-level and dog-
263 level effects. The mean model, describing the predicted pattern of behaviour across days
264 on the latent scale, y^* , for observation i from dog j , was modelled as:

$$y_{ij}^* = \beta_0 + v_{0j} + \sum_{p=1}^P \beta_{p0} x_{pj} + \left(\beta_1 + v_{1j} + \sum_{p=1}^P \beta_{p1} x_{pj} \right) day_{ij} + \left(\beta_2 + v_{2j} + \sum_{p=1}^P \beta_{p2} x_{pj} \right) day_{ij}^2 \quad (2)$$

265 Equation 2 expresses the longitudinal pattern of behaviour as a function of i) a group-
266 level intercept the same for all dogs, β_0 , and the deviation from the group-level intercept
267 for each dog, v_{0j} , ii) a linear effect of day since arrival, β_1 , and each dog's deviation, v_{1j} ,
268 and iii) a quadratic effect of day since arrival, β_2 , and each dog's deviation, v_{2j} . A
269 quadratic effect was chosen based on preliminary plots of the data at group-level and at
270 the individual-level, although we also compared the model's predictive accuracy with
271 simpler models (described below). Day since arrival was standardised, meaning that the
272 intercepts reflected the behaviour on the average day since arrival across dogs
273 (approximately day 8). The three dog-level parameters, v_j , correspond to personality and
274 linear and quadratic plasticity parameters, respectively. The terms $\sum_{p=1}^P \beta_p x_{pj}$ denote the
275 effect of P dog-level predictor variables (x_p), included to explain variance between dog-
276 level intercepts and slopes. These included: the number of observations for each dog, the
277 number of days dogs spent at the shelter controlling for the number of observations (i.e.
278 the residuals from a linear regression of total number of days spent at the shelter on the
279 number of observations), average age while at the shelter, average weight at the shelter,
280 sex, neuter status, source type, and rehoming centre (Table 1). For neuter status, we did
281 not make comparisons between the 'undetermined' category and other categories. The

282 primary goal of including these predictor variables was to obtain estimates of individual
 283 differences conditional on relevant inter-individual differences variables, since the data
 284 were observational.

285

286 The SD model was:

$$\sigma = \exp \left(\delta + v_{3j} + \sum_{p=1}^P \beta_{p3} x_{pj} \right) \quad (3)$$

287 This equation models the SD of the latent scale by its own regression, with group-level
 288 SD intercept, δ , evaluated at the average day, the deviation for each dog from the group-
 289 level SD intercept, v_{3j} , and predictor variables, $\sum_{p=1}^P \beta_{p3} x_{pj}$, as in the mean model
 290 (equation 2). The SDs across dogs were assumed to approximately follow a log-normal
 291 distribution, with $\ln(\sigma)$ approximately normally distributed (hence the exponential
 292 inverse-link function). The parameter v_{3j} corresponds to each dog's residual SD or
 293 predictability.

294

295 All four dog-level parameters were assumed to be multivariate normally distributed with
 296 means 0 and variance-covariance matrix $\Sigma_{\mathbf{v}}$ estimated from the data:

$$\Sigma_{\mathbf{v}} = \begin{bmatrix} \tau_{v_0}^2 & \rho_{v_{01}} \tau_{v_0} \tau_{v_1} & \rho_{v_{02}} \tau_{v_0} \tau_{v_2} & \rho_{v_{03}} \tau_{v_0} \tau_{v_3} \\ \dots & \tau_{v_1}^2 & \rho_{v_{12}} \tau_{v_1} \tau_{v_2} & \rho_{v_{13}} \tau_{v_1} \tau_{v_3} \\ \dots & \dots & \tau_{v_2}^2 & \rho_{v_{23}} \tau_{v_2} \tau_{v_3} \\ \dots & \dots & \dots & \tau_{v_3}^2 \end{bmatrix} \quad (4)$$

297 The diagonal elements are the variances of the dog-level intercepts, linear slopes,
 298 quadratic slopes and residual SDs, respectively, while the covariances fill the off-
 299 diagonal elements (only the upper triangle shown), where ρ is the correlation coefficient.
 300 In the results, we report τ_{v_3} (the SD of dog-level residual SDs) on the original scale,

301 rather than the log-transformed scale, using $\sqrt{e^{2\delta + \tau_{v_3}^2} e^{\tau_{v_3}^2} - 1}$. Likewise, δ was
302 transformed to the median of the original scale by e^δ .

303

304 To summarise the amount of behavioural variation explained by differences between
305 individuals, referred to as repeatability in the personality literature [1], we calculated the
306 intra-class correlation coefficient (ICC). Since the model includes both intercepts and
307 slopes varying by dog, the ICC is a function of both linear and quadratic effects of day
308 since arrival. The ICC for day i , assuming individuals with the same residual variance
309 (i.e. using the median of the log-normal residual SD), was calculated as:

$$ICC_i = \frac{\tau_{v_0}^2 + 2Cov_{v_0, v_1} Day_i + \tau_{v_1}^2 Day_i^2 + 2Cov_{v_0, v_2} Day_i^2 + \tau_{v_2}^2 Day_i^4 + 2Cov_{v_1, v_2} Day_i^3}{numerator + e^\delta} \quad (5)$$

310 Equation 5 is an extension of the intra-class correlation calculated from mixed-effect
311 models with a random intercept only [43] to include the variance parameters for, and
312 covariances between, the linear and quadratic effects of day, which were evaluated at
313 specific days of interest. We calculated the ICC for values of -1, 0 and 1 on the
314 standardised day scale, corresponding to approximately the arrival day (day 0), day 8, and
315 day 15. This provided a representative spread of days for most of the dogs in the sample,
316 since there were fewer data available for later days which could lead to inflation of inter-
317 individual differences.

318

319 To inspect the degree of rank-order change in sociability across dogs from arrival day
320 compared to specific later days (i.e. whether dogs that were, on average, least sociable on
321 arrival also tended to be least sociable later on), we calculated the ‘cross-environmental’
322 correlations [44] between the same days as the ICC. The cross-environmental covariance
323 matrix, Ω , between the three focal days was calculated as:

$$\Omega = \Psi K \Psi' \quad (6)$$

324

325 In equation 6, \mathbf{K} represents the variance-covariance matrix of the dog-level intercepts and
326 (linear and quadratic) slopes, and $\mathbf{\Psi}$ is a three-by-three matrix with a column vector of 1s,
327 a column vector containing -1, 0, and 1 defining the day values for the cross-
328 environmental correlations for the linear component, and a column vector containing 1, 0,
329 and 1 defining the day values for the cross-environmental correlations for the quadratic
330 component. Once defined, $\mathbf{\Omega}$ was scaled to a correlation matrix. Finally, to summarise the
331 degree of individual differences in predictability, we calculated the ‘coefficient of
332 variation for predictability’ as $\sqrt{e^{\tau^2 v_3} - 1}$ following Cleasby *et al.* [5].

333

334 **Prior distributions**

335 We chose prior distributions that were either weakly informative (i.e. specified a realistic
336 range of parameter values) for computational efficiency, or weakly regularising to
337 prioritise conservative inference. The prior for the overall intercept, β_0 , was
338 $Normal(\bar{y}, 5)$, where \bar{y} is the arithmetic mean of the ordinal data. The linear and
339 quadratic slope parameters, β_1 and β_2 , were given $Normal(0, 1)$ priors. Coefficients for
340 the dog-level predictor variables, β_k , were given $Normal(0, \sigma_{\beta_p})$ priors, where σ_{β_p} was a
341 shared SD across predictor variables, which had in turn a half-Cauchy hyperprior with
342 mode 0 and shape parameter 2, $half - Cauchy(0, 2)$. Using a shared SD imposes
343 shrinkage on the regression coefficients for conservative inference: when most regression
344 coefficients are near zero, then estimates for other regression coefficients are also pulled
345 towards zero (e.g. [32]). The prior for the overall log-transformed residual SD, δ , was
346 $Normal(0, 1)$. The covariance matrix of the random effects was parameterised as a
347 Cholesky decomposition of the correlation matrix (see [45] for more details), where the
348 SDs had $half - Cauchy(0, 2)$ priors and the correlation matrix had a LKJ prior
349 distribution [46] with shape parameter η set to 2.

350

351 **Model selection & computation**

352 We compared the full model explained above to five simpler models. Starting with the
353 full model, the alternative models included: i) parameters quantifying personality and
354 quadratic and linear plasticity only; ii) parameters quantifying personality and linear
355 plasticity only, with a fixed quadratic effect of day since arrival; iii) parameters
356 quantifying personality only, with fixed linear and quadratic effects of day since arrival;
357 iv) parameters quantifying personality only, with a fixed linear effect of day since arrival;
358 and v) a generalised linear regression with no dog-varying parameters and a linear fixed
359 effect for day since arrival (Figure 1). Models were compared by calculating the widely
360 applicable information criterion (WAIC; [47]) following McElreath [33] (see the R script
361 file). The WAIC is a fully Bayesian information criterion that indicates a model's *out-of-*
362 *sample* predictive accuracy relative to other plausible models while accounting for model
363 complexity, and is preferable to the deviance information criterion (DIC) because WAIC
364 does not assume multivariate normality in the posterior distribution and returns a
365 probability distribution rather than a point estimate [33]. Thus, WAIC guards against both
366 under- and over-fitting to the data (unlike measures of purely in-sample fit, e.g. R^2).

367

368 Models were computed using the probabilistic programming language Stan [45] using the
369 *RStan* package [48] version 2.15.1, which employs Markov chain Monte Carlo estimation
370 using Hamiltonian Monte Carlo (see the R script file and Stan code for full details). We
371 ran four chains of 5,000 iterations each, discarding the first 2,500 iterations of each chain
372 as warm-up, and setting thinning to 1. Convergence was assessed visually using trace
373 plots to ensure chains were well mixed, numerically using the Gelman-Rubin statistic
374 (values close to 1 and < 1.05 indicating convergence) and by inspecting the effective
375 sample size of each parameter. We also used graphical posterior predictive checks to
376 assess model predictions against the raw data, including 'counterfactual' predictions [33]
377 to inspect how dogs would be predicted to behave across the first month of being in the
378 shelter regardless of their actual number of observations or length of stay at the shelter.

379 To summarise parameter values, we calculated mean (denoted β) and 95% highest
380 density intervals (HDIs), the 95% most probable values for each parameter (using
381 functions in the *rethinking* package; [33]). For comparing levels of categorical variables,
382 the 95% HDI of their differences were calculated (i.e. the differences between the
383 coefficients at each step in the MCMC chain, denoted β_{diff}). When the 95% HDI of
384 predictor variables surpassed zero, a credible effect was inferred.

385

386 **Results**

387 **Inter-rater reliability & validity**

388 For the two videos depicting interactions with people, consensus was 0.75 (95% CI: 0.66,
389 0.84) for the video showing an example of *Reacts to people non-aggressive* and 0.77
390 (95% CI: 0.74, 0.81) for the example of *Reacts to people aggressive*, respectively.

391 Neither did these results overlap with the null distributions (see Supplementary Material
392 Table S1), indicating significant inter-rater reliability. For the video showing *Reacts to*
393 *people non-aggressive*, 77% chose the correct code and 83% a code of the correct colour
394 category (amber), and, as previously reported by [35], 52% chose the correct code for the
395 video showing *Reacts to people aggressive* and 55% chose a code of the correct colour
396 category (red; 42% chose the amber code *Reacts to people non-aggressive* instead).

397 Across all assessment context videos, the average consensus was 0.71 and participants
398 chose the correct ethogram category 66% of the time while 78% of answers were a
399 category of the correct ethogram colour.

400

401 **Hierarchical ordinal probit model**

402 The full model had the best out-of-sample predictive accuracy, with the inclusion of
403 heterogeneous residual SDs among dogs improving model fit by over 1,500 WAIC points

404 compared to the second most plausible model (Alternative 1 in Figure 1). In general,
405 models that included more parameters to describe personality, plasticity and
406 predictability, and models with a quadratic effect of day, had better out-of-sample
407 predictive accuracy, despite the added complexity brought by additional parameters.

408

409 At the group-level, the *Friendly* code (Table 2) was most probable overall and was
410 estimated to increase in probability across days since arrival, while the remaining
411 sociability codes either decreased or stayed at low probabilities (Figure 2a), reflecting the
412 raw data. On the latent sociability scale (Figure 2b), the group-level intercept parameter
413 on the average day was 0.68 (95% HDI: 0.51, 0.86). A one SD increase in the number of
414 days since arrival was associated with a -0.63 unit (95% HDI: -0.77, -0.50) change on the
415 latent scale on average (i.e. reflecting increasing sociability), and the group-level
416 quadratic slope was positive ($\beta = 0.20$, 95% HDI: 0.10, 0.30), reflecting a quicker rate of
417 change in sociability earlier after arrival to the shelter than later (i.e. a concave down
418 parabola). There was a slight increase in the quadratic curve towards the end of the one-
419 month period, although there were fewer behavioural observations at this point and so
420 greater uncertainty about the exact shape of the curve, resulting in estimates being pulled
421 closer to those of the intercepts. The group-level residual standard deviation had a median
422 of 1.84 (95% HDI: 1.67, 2.02).

423

424 At the individual level, heterogeneity existed in behavioural trajectories across days since
425 arrival (Figure 2b). The SDs of dog-varying parameters were: i) intercepts: 1.29 (95%
426 HDI: 1.18, 1.41; Figure 3a), ii) linear slopes: 0.56 (95% HDI: 0.47, 0.65; Figure 3b), iii)
427 quadratic slopes: 0.28 (95% HDI: 0.20, 0.35; Figure 3c), and iv) residual SDs: 1.39 (95%
428 HDI: 1.22, 1.58; Figure 3d). There was also large uncertainty in individual-level
429 estimates. Figure 4 displays counterfactual model predictions for twenty randomly-
430 sampled dogs. Uncertainty in reaction norm estimates, illustrated by the width of the 95%
431 HDIs (dashed black lines), was greatest when data were sparse (e.g. towards the end of
432 the one-month study period). Hierarchical shrinkage meant that individuals with

433 observations of less sociable responses, or individuals with few behavioural observations,
434 tended to have model predictions pulled towards the overall mean. Note that regression
435 lines depict values on the latent scale predicted to generate observations on the ordinal
436 scale, and so may not clearly fit the ordinal data points. The coefficient of variation for
437 predictability was 0.64 (95% HDI: 0.58, 0.70). Individuals with the five highest and
438 lowest residual SD estimates are shown in Figure 5.

439

440 Dog-varying intercepts positively correlated with linear slope parameters ($\rho = 0.38$, 95%
441 HDI: 0.24, 0.50) and negatively correlated with quadratic slope parameters ($\rho = -0.54$,
442 95% HDI: -0.68, -0.39), and linear and quadratic slopes had a negative correlation ($\rho = -$
443 0.75 , 95% HDI: -0.88, -0.59), indicating that less sociable individuals (with higher scores
444 on the ordinal scale) had flatter reaction norms on average. Dog-varying residual SDs had
445 a correlation with the intercept parameters of approximately zero ($\rho = 0.00$, 95% HDI: -
446 0.10 , 0.10) but were negatively correlated with the linear slope parameters ($\rho = -0.37$,
447 95% HDI: -0.51, -0.22) and positively correlated with the quadratic slopes ($\rho = 0.24$,
448 95% HDI: 0.05, 0.42), indicating that dogs with greater residual SDs were predicted to
449 change the most across days since arrival.

450

451 The ICC by day increased from arrival day (ICC = 0.22; 95% HDI: 0.16, 0.28) to day 8
452 (ICC = 0.33; 95% HDI: 0.28, 0.38) but changed little by day 15 (ICC = 0.32; 95% HDI:
453 0.27, 0.37). The cross-environmental correlation between days 0 and 8 was 0.79 (95%
454 HDI: 0.70, 0.88), between days 0 and 15 was 0.51 (95% HDI: 0.35, 0.68), and between
455 days 8 and 15 was 0.95 (95% HDI: 0.93, 0.97).

456

457 A one SD increase in the number of observations was associated with higher intercepts
458 ($\beta = 0.12$; 95% HDI: 0.03, 0.21; see Supplementary Material Table S2) and higher
459 residual SDs ($\beta = 0.06$, 95% HDI: 0.02, 0.10). Increasing age by one SD was associated
460 with lower intercepts ($\beta = -0.61$, 95% HDI: -0.70, -0.51), steeper linear slopes ($\beta = -0.20$,

461 95% HDI: -0.27, -0.13), a stronger quadratic curve ($\beta = 0.07$, 95% HDI: 0.03, 0.12), and
462 larger residual SDs ($\beta = 0.05$, 95% HDI: 0.01, 0.09). Increasing weight by one SD was
463 associated with shallower quadratic curves ($\beta = -0.05$, 95% HDI: -0.09, -0.01). No
464 credible effect of sex was observed on personality, plasticity or predictability. Gift dogs
465 had larger intercepts than returned dogs ($\beta_{diff} = 0.28$, 95% HDI: 0.04, 0.52) and stray
466 dogs ($\beta_{diff} = 0.33$, 95% HDI: 0.15, 0.50), as well as steeper linear slopes ($\beta_{diff} = -0.25$,
467 95% HDI: -0.38, -0.13) and higher residual SDs than stray dogs ($\beta_{diff} = 0.10$, 95% HDI:
468 0.02, 0.18). Dogs at the large rehoming centre had steeper linear slopes ($\beta_{diff} = -0.70$,
469 95% HDI: -0.84, -0.56) and stronger quadratic curves ($\beta_{diff} = 0.35$, 95% HDI: 0.26,
470 0.45) than dogs at the medium rehoming centre, and lower intercept parameters ($\beta_{diff} = -$
471 0.30, 95% HDI: -0.50, -0.09) and steeper linear slopes ($\beta_{diff} = -0.22$, 95% HDI: -0.38, -
472 0.06) than dogs at the small rehoming centre. Compared to dogs at the small rehoming
473 centre, dogs at the medium centre had lower intercepts ($\beta_{diff} = -0.25$, 95% HDI: -0.48, -
474 0.01), and shallower linear ($\beta_{diff} = 0.48$, 95% HDI: 0.30, 0.66) and quadratic slopes
475 ($\beta_{diff} = -0.34$, 95% HDI: -0.46, -0.22). Dogs already neutered before arrival to the
476 shelter had lower intercepts ($\beta_{diff} = -0.54$, 95% HDI: -1.07, -0.03) and lower residual
477 SDs ($\beta_{diff} = -0.53$, 95% HDI: -0.85, -0.22) than dogs not neutered, but higher intercepts
478 ($\beta_{diff} = 0.20$, 95% HDI: 0.03, 0.37) and higher residual SDs ($\beta_{diff} = 0.10$, 95% HDI:
479 0.02, 0.19) than those neutered whilst at the shelter. Unneutered dogs had higher
480 intercepts ($\beta_{diff} = 0.74$, 95% HDI: 0.20, 1.26) and higher residual SDs ($\beta_{diff} = 0.63$,
481 95% HDI: 0.30, 0.92) than dogs neutered at the shelter.

482

483 **Discussion**

484 This study applied the framework of behavioural reaction norms to quantify inter- and
485 intra-individual differences in shelter dog behaviour during interactions with unfamiliar
486 people. This is the first study to systematically analyse behavioural data from a

487 longitudinal, observational assessment of shelter dogs. Dogs demonstrated substantial
488 individual differences in personality, plasticity and predictability, which were not well
489 described by simply investigating how dogs behaved on average. In particular,
490 accounting for individual differences in predictability, or the short-term, day-to-day
491 fluctuations in behaviour, resulted in significant improvement in model fit (Figure 1).
492 Modelling dogs' longitudinal behaviour also demonstrated that behavioural repeatability
493 increased with days since arrival (i.e. increasing proportion of variance explained by
494 between-individual differences), particularly across the first week since arrival. Similarly,
495 while individuals maintained rank-order differences in sociability across smaller periods
496 (i.e. first 8 days), rank-order differences were only moderately maintained between
497 arrival at the shelter and day 15. The results highlight the importance of adopting
498 observational and longitudinal assessments of shelter dog behaviour, provide a method by
499 which to analyse longitudinal data commensurate with other work in animal behaviour,
500 and identify previously unconsidered behavioural measures that could be used to improve
501 the predictive validity of behavioural assessments in dogs.

502

503 **Average behaviour**

504 At the group-level, dogs' reactions to meeting unfamiliar people were predominantly
505 coded as *Friendly* (Figure 2a), described as 'Dog initiates interactions in an appropriate
506 social manner'. Although this definition is broad, it represents a functional qualitative
507 characterisation of behaviour suitable for the purposes of the shelter when coding
508 behavioural interactions, and its generality may partly explain why it was the most
509 prevalent category. The results are consistent with findings that behaviours indicative of
510 poor welfare and/or difficulty of coping (e.g. aggression) are relatively infrequent even in
511 the shelter environment [22,26]. The change of behaviour across days since arrival was
512 characterised by an increase in the *Friendly* code and a decrease in other behavioural
513 codes (Figure 2a). Furthermore, the positive quadratic effect of day since arrival on
514 sociability illustrates that the rate of behavioural change was not constant across days,
515 being quickest earlier after arrival (Figure 2b). The range of behavioural change at the

516 group-level was, nevertheless, still concentrated around the lowest behavioural codes,
517 *Friendly* and *Excitable*.

518

519 Previous studies provide conflicting evidence regarding how shelter dogs adapt to the
520 kennel environment over time, including behavioural and physiological profiles
521 indicative of both positive and negative welfare [26]. Whereas some authors report
522 decreases in the prevalence of some stress- and/or fear related behaviour with time
523 [27,49], others have reported either no change or an increase in behaviours indicative of
524 poor welfare [17,30]. Of relevance here, Kis *et al.* [17] found that aggression towards
525 unknown people increased over the first two weeks of being at a shelter. In the current
526 study, aggression was rare (Table 2), and the probability of ‘red codes’ (which included
527 aggression) decreased with days at the shelter (Figure 3a). A salient difference is that Kis
528 *et al.* [17] collected data using a standardised behavioural test consisting of a stranger
529 engaging in a ‘threatening approach’ towards dogs. By contrast, we used a large data set
530 of behavioural observations recorded after non-standardised, spontaneous interactions
531 between dogs and unfamiliar people. In recording spontaneous interactions, the shelter
532 aimed to elicit behaviour more representative of a dog’s typical behaviour outside of the
533 shelter environment than would be seen in a standardised behavioural assessment.
534 Previously, authors have noted that standardised behavioural assessments may induce
535 stress and inflate the chances of dogs displaying aggression [29], emphasising the value
536 of observational methods of assessment in shelters [24]. While such observational
537 methods are less standardised, they may have greater ecological validity by giving results
538 more representative of how dogs will behave outside of the shelter. Testing the predictive
539 value of observational assessments on behaviour post-adoption is the focus of ongoing
540 research.

541

542 **Individual-level variation**

543 When behavioural data are aggregated across individuals, results may provide a poor
544 representation of how individuals in a sample actually behaved. Here, we found

545 heterogeneity in dog behaviour across days since arrival, even after taking into account a
546 number of dog-level predictor variables that could explain inter-individual differences.
547 Variation in individuals' average behaviour across days (i.e. variation in dogs' intercept
548 estimates) illustrated that personality estimates spanned a range of behavioural codes,
549 although model predictions mostly spanned the green codes (Figure 2b; Table 2).
550 However, whilst there were many records to inform group-level estimates, there were
551 considerably fewer records available for each individual, which resulted in large
552 uncertainty of individual personality parameters (illustrated by wide 95% HDI bars in
553 Figure 3a). Personality variation has been the primary focus of previous analyses of
554 individual differences in dogs, often based on data collected at one time point and usually
555 on a large number of behavioural variables consolidated into composite or latent
556 variables (e.g. [50–52]). Our results highlight that ranking individuals on personality
557 dimensions from few observations entails substantial uncertainty.

558

559 Certain studies on dog personality have explored how personality trait scores change
560 across time periods, such as ontogeny (e.g. [53]) or time at a shelter (e.g. [17]). Such
561 analyses assume, however, that individuals have similar degrees of change through time.
562 If individuals differ in the magnitude or direction of change (i.e. degree of plasticity),
563 group-level patterns of change may not capture important individual heterogeneity. In
564 this study, most dogs were likely to show lower behavioural codes/more sociable
565 responses across days since arrival, although the rate of linear and quadratic change
566 differed among dogs. Indeed, some dogs showed a *decrease* in sociability through time
567 (individuals with positive model estimates in Figure 3b), and while most dogs showed
568 greater behavioural change early after arrival, others showed slower behavioural change
569 early after arrival (individuals with negative model estimates in Figure 3c). As with
570 estimates of personality, there was also large uncertainty of plasticity.

571

572 Part of the difficulty of estimating reaction norms for heterogeneous data is choosing a
573 function that best describes behavioural change. We examined both linear and quadratic

574 effects of day since arrival based on preliminary plots of the data, and their inclusion in
575 the best fitting full model is supported by the lower WAIC value of alternative model 3,
576 with both effects, compared to 4, with just the linear effect (Figure 1). Most studies are
577 constrained to first-order polynomial reaction norms through time due to collecting data
578 at only a few time points [6,44]. However, the quadratic function was relatively easy to
579 vary across individuals while maintaining interpretability of the results. More complex
580 functions (e.g. regression splines) have the disadvantage of being less easily interpretable
581 and higher-order polynomial functions may produce only crude representations of data-
582 generating processes [33]. Nevertheless, by collecting data more intensely, the
583 opportunities to model behavioural reaction norms beyond simple polynomial effects of
584 time should improve. For instance, ecological momentary assessment studies in
585 psychology point to possibilities for modelling behaviour as a dynamic system, such as
586 with the use of vector-autoregressive models and dynamic network or factor models (e.g.
587 [54,55]). These models can also account for relationships between multiple dependent
588 variables (e.g. multiple measures of sociability). Models of behavioural reaction norms,
589 by contrast, have usually been applied to only one dependent variable operationally
590 defined as reflecting the trait of interest, so methods to model multiple dependent
591 variables through time concurrently will be an important advancement.

592

593 Personality and plasticity were correlated, with dogs with less sociable behaviour across
594 days being less plastic. Previous studies have explored the relationship between how
595 individuals behave on average and their degree of behavioural change. David *et al.* [56]
596 found that male golden hamsters (*Mesocricetus auratus*) showing high levels of
597 aggression in a social intruder paradigm were slower in adapting to a delayed-reward
598 paradigm. In practice, the relationship between personality and plasticity is probably
599 context dependent. Betini and Norris [57] found, for instance, that more aggressive male
600 tree swallows (*Tachycineta bicolor*) during nest defence were more plastic in response to
601 variation in temperature, but that plasticity was only advantageous for nonaggressive
602 males and no relationship was present between personality and plasticity in females. The
603 correlation between personality and plasticity indicates a ‘fanning out’ shape of the

604 reaction norms through time (Figure 2b). Consequently, behavioural repeatability or the
605 amount of variance explained by between-individual differences increased as a function
606 of day, but only after the first week after arrival. The ‘cross-environmental’ correlation,
607 moreover, indicated that the most sociable dogs on arrival day were not necessarily the
608 most sociable on later days at the shelter. In particular, the correlation between sociability
609 scores on arrival day and day 15 was only moderate, supporting Brommer [44] that the
610 rank-ordering of trait scores is not always reliable. By contrast, the cross-environmental
611 correlations between days 0 and 8, and between days 8 and 15, were much stronger.
612 These results suggest that shelters using standardised behavioural assessments would
613 benefit from administering such tests as late as possible after dogs arrive.

614

615 Of particular interest was predictability or the variation in dogs’ residual SDs. Studies of
616 dog personality generally treat behaviour as probabilistic, implying recognition that
617 residual intra-individual behaviour is not completely stable, and authors have posited that
618 dogs may vary in their behavioural consistency (e.g. [13]). Yet, this is the first study to
619 quantify individual differences in predictability in dogs. Modelling residual SDs for each
620 dog resulted in a model with markedly better out-of-sample predictive accuracy (Figure
621 1). The coefficient of variation for predictability was 0.64 (95% HDI: 0.58, 0.70), which
622 is high compared to other studies in animal behaviour. For instance, Mitchell *et al.* [6]
623 reported a value of 0.43 (95% HDI: 0.36, 0.53) in spontaneous activity measurements of
624 male guppies (*Poecilia reticulata*). Variation in predictability also supports the
625 hypothesis that dogs have varying levels of behavioural consistency. It is important to
626 note, however, that interactions with unfamiliar people at the shelter were likely more
627 heterogeneous than behavioural measures from standardised tests or laboratory
628 environments, which may contribute to greater individual variation in predictability.
629 Moreover, the behavioural data analysed here may have contained more measurement
630 error than data from more standardised environments.

631

632 Although shelter employees demonstrated significant inter-rater reliability in video
633 coding sessions, the average proportion of shelter employees who selected the correct
634 behavioural code to describe behaviours seen in videos was modest (66%), while 78%
635 chose a video in the correct colour category (green, amber or red). Indeed, only 55% of
636 employees identified the *Reacts to people aggressive* behaviour as a red code, with the
637 remaining employees identifying it as the amber category code *Reacts to people non-*
638 *aggressive*. As discussed by Goold and Newberry [35], employees were likely to mistake
639 examples of aggression for non-aggression, but not the other way around. In the current
640 study, this would have increased the percentage of lower category codes (describing
641 greater sociability). Due to lower standardisation of the observational contexts at the
642 shelter than in formal behavioural testing, it was important to evaluate the reliability and
643 validity of the behavioural records. Defining acceptable standards of reliability and
644 validity is, however, non-trivial and we could not find measures of reliability or validity
645 in any previous studies investigating predictability in animals for comparison.

646

647 Dogs with higher residual SDs demonstrated steeper linear slopes and greater quadratic
648 curves, indicating that greater plasticity was associated with lower predictability. The
649 costs of plasticity are believed to include greater phenotypic instability, in particular
650 developmental instability [11,58]. Since more plastic individuals are more responsive to
651 environmental perturbation, a limitation of plasticity may be greater phenotypic
652 fluctuation on finer time scales. However, lower predictability may also confer a benefit
653 to individuals precisely because they are less predictable to con- and hetero-specifics. For
654 instance, Highcock and Carter [59] reported that predictability in behaviour decreases
655 under predation risk in Namibian rock agamas (*Agama planiceps*). No correlation was
656 found here between personality and predictability, similar to findings of Biro and
657 Adriaenssens [2] in mosquitofish (*Gambusia holbrooki*), although correlations were
658 found in agamas [59] and guppies [6]. It is possible that correlations between personality
659 and predictability depend upon the specific aspects of personality under investigation.

660

661 **Predictors of individual variation**

662 Finally, we found associations between certain predictor variables and personality,
663 plasticity and predictability (Supplementary Material Table S2). Our primary reason for
664 including these predictor variables was to obtain more accurate estimates of personality,
665 plasticity and predictability, and we remain cautious about *a posteriori* interpretations of
666 their effects, especially since the theory underlying why individuals may, for example,
667 demonstrate differences in predictability is in its infancy [8]. The reproducibility of a
668 number of the results would, nevertheless, be interesting to confirm in future research. In
669 particular, understanding factors affecting intra-individual change is important given that
670 many personality assessments are used to predict an individual's future behaviour, rather
671 than understand inter-individual differences. Here, increasing age was associated with
672 greater plasticity (linear and quadratic change) and lower predictability, although some of
673 the parameters' 95% HDIs were close to zero, indicative of small effects. In great tits
674 (*Parus major*) conversely, plasticity decreased with age [60], whilst in humans, intra-
675 individual variability in reaction times increased with age [61]. Moreover, non-neutered
676 dogs showed lower predictability than neutered dogs, and dogs entering the shelter as
677 gifts (relinquished by their owners) had lower predictability estimates than stray dogs
678 (dogs brought in by local authorities or members of the public after being found without
679 their owners). These results can be used to formulate specific hypotheses about
680 behavioural variation.

681 **Conclusion**

682 We applied the framework of behavioural reactions norms to data from a longitudinal and
683 observational shelter dog behavioural assessment, quantifying inter- and intra-individual
684 behavioural variation in dogs' interactions with unfamiliar people. Overall, shelter dogs
685 were sociable with unfamiliar people and sociability continued to increase with days
686 since arrival to the shelter. At the same time, dogs showed individual differences in
687 personality, plasticity and predictability. Accounting for all of these components
688 substantially improved model fit, particularly the inclusion of predictability, which

689 suggests that individual differences in day-to-day behavioural variation represent an
690 important, yet largely unstudied, component of dog behaviour. Our results also highlight
691 the uncertainty of making predictions about shelter dog behaviour, particularly when the
692 number of behavioural observations is low. For shelters conducting standardised
693 behavioural assessments, assessments are likely best carried out as late as possible, given
694 that rank-order differences between individuals on arrival and at day 15 were only
695 moderately related. In conclusion, this study supports moving towards observational and
696 longitudinal assessments of shelter dog behaviour, has demonstrated a Bayesian method
697 by which to analyse longitudinal data on dog behaviour, and suggests that the predictive
698 validity of behavioural assessments in dogs could be improved by systematically
699 accounting for both inter- and intra-individual variation.

700 **Ethics statement**

701 Full permission to use the data in this article was provided by Battersea Dogs and Cats
702 Home.

703 **Data accessibility**

704 The data, R code and Stan model code to run the analyses and produce the results and
705 figures in this article are available on Github:
706 https://github.com/ConorGoold/GooldNewberry_modelling_shelter_dog_behaviour

707 **Competing interests**

708 We declare no competing interests.

709 **Author contributions**

710 CG and RCN conceptualised the study. CG obtained the data, conducted the statistical
711 analyses and drafted the initial manuscript. CG and RCN revised the manuscript and
712 wrote the final version.

713 **Acknowledgements**

714 The authors are especially grateful to Battersea Dogs and Cats Home for providing the
715 data on their behavioural assessment.

716 **Funding statement**

717 CG and RCN are employed by the Norwegian University of Life Sciences. No additional
718 funding was required for this study.

719

720

721 **Table 1.** Demographic variables of dogs in the sample analysed. Mean and standard
722 deviation (SD) or the number of dogs by category (N) are displayed.

723 **Table 2.** Ethogram of behavioural codes used to record observations of interactions with
724 unfamiliar people, and their percent prevalence in the sample. Behaviour labels followed
725 by + indicate a more intense form of the behaviour with the same name without a +.

726 **Figure 1.** Out-of-sample predictive accuracy (lower is better) for each model (described
727 in text section section 2.5.5) measured by the widely applicable information criterion
728 (WAIC). Black points denote the WAIC estimate and horizontal lines show WAIC
729 estimates \pm standard error. Mean \pm standard error: full model = 38669 ± 275 ; alternative
730 1 = 40326 ± 288 ; alternative 2 = 40621 ± 288 ; alternative 3 = 40963 ± 289 ; alternative 4
731 = 41100 ± 289 ; alternative 5 = 45268 ± 289 .

732 **Figure 2.** (a) Predicted probabilities (posterior means = black lines; 95% highest density
733 intervals = shaded areas) of different sociability codes across days since arrival. (b)
734 Posterior mean behavioural trajectories on the latent scale (ranging from $\pm\infty$) at the
735 group-level (blue line) and for each individual (black lines), where higher values indicate
736 lower sociability.

737 **Figure 3.** Posterior means (black dots) and 95% highest density intervals (grey vertical
738 lines) for each dogs' (a) intercept, (b) linear slope, (c) quadratic slope, and (d) residual
739 SD parameter.

740 **Figure 4.** Predicted reaction norms ('counterfactual' plots) for twenty randomly-selected
741 dogs. Black points show raw data on the ordinal scale (higher values indicate lower
742 sociability), and solid and dashed lines illustrate posterior means and 95% highest density
743 intervals. When data were sparse, there was increased uncertainty in model predictions.
744 Due to hierarchical shrinkage, individual dogs' model predictions were pulled towards
745 the group-level mean, particularly for those dogs showing higher behavioural codes (i.e.
746 less sociable responses).

747 **Figure 5.** Reaction norms (posterior means = solid black lines; 95% highest density
748 intervals = dashed black lines) for individuals with the five highest (top row) and five
749 lowest (bottom row) residual SDs. Black points represent raw data on the ordinal scale
750 (higher values indicating lower sociability).

751

Demographic variable	Mean (SD) / N
Number of observations per dog	5.9 (3.7)
Days spent at the shelter	25.8 (35.0)
Age (years; all at least 4 months old)	3.7 (3.0)
Weight (kg)	18.9 (10.2)
Source: gift / stray / return	1950 / 1122 / 191
Rehoming centre: London / Old Windsor / Brands Hatch	1873 / 951 / 439
Females / males	1396 / 1867
Neutered: before arrival / at shelter / not / undetermined	1043 / 1281 / 747 / 192

752

753

Behaviour	Colour	%	Definition
1: Friendly	Green	63.5	Dog initiates interactions with people in an appropriate social manner.
2: Excitable	Green	14.2	Animated interaction with an enthusiastic attitude, showing behaviours such as jumping up, mouthing, an inability to stand still, and/or playful behaviour towards people.
3: Independent	Green	4.1	Does not actively seek interaction, although relaxed in the presence of people
4: Submissive	Green	4.6	Appeasing and/or nervous behaviours, including a low body posture, rolling over and other calming signals.
5: Uncomfortable avoids	Amber	5.4	Tense and stiff posture, and/or shows anxious behaviours (e.g. displacement behaviours) while trying to move away from the person.
6: Submissive +	Amber	0.2	High intensity of submissive behaviours such as submissive urination, a reluctance to move, or is frequently overwhelmed by the interaction.
7: Uncomfortable static	Amber	0.8	Tense and stiff posture, and/or shows anxious behaviour (potentially showing displacement behaviours) but doesn't move away from the person.
8: Stressed	Amber	0.5	High frequency/intensity of stress behaviours, which may include dribbling, stereotypic behaviours, stress vocalisations, constant shedding, trembling, and destructive behaviours.
9: Reacts to people non-aggressive	Amber	2.4	Barks, whines, howls and/or play growls when seeing/meeting people, potentially pulling or lunging towards them.
10: Uncomfortable approaches	Amber	0.7	Tense and stiff posture, and/or shows anxious behaviour (potentially showing displacement behaviours) and approaches the person.
11: Overstimulated	Red	0.8	High intensity of excitable behaviour, including grabbing, body barging, and nipping.
12: Uncomfortable static +	Red	0.1	Body freezes (the body goes suddenly and completely still) in response to an interaction with a person.
13: Reacts to people aggressive	Red	2.8	Growls, snarls, shows teeth and/or snaps when seeing/meeting people, potentially pulling or lunging towards them.

754

755

- 756 1. Bell, A. M., Hankison, S. J. & Laskowski, K. L. 2009. The repeatability of behaviour:
757 a meta-analysis. *Anim. Behav.* **77**, 771–783. (doi: 10.1016/j.anbehav.2008.12.022.□)
- 758 2. Biro, P. A., Adriaenssens, B., Cole, A. E. B. J. & Bronstein, E. J. L. 2013.
759 Predictability as a personality trait: consistent differences in intraindividual behavioral
760 variation. *Am. Nat.* **182**, 621–629. (doi:10.1086/673213)
- 761 3. Dingemanse, N. J., Kazem, A. J. N., Réale, D. & Wright, J. 2010. Behavioural reaction
762 norms: animal personality meets individual plasticity. *Trends Ecol. Evol.* **25**, 81–89.
763 (doi:10.1016/j.tree.2009.07.013)
- 764 4. Bridger, D., Bonner, S. J. & Briffa, M. 2015 Individual quality and personality: bolder
765 males are less fecund in the hermit crab (*Pagurus bernhardus*). *Proc. R. Soc. B* **282**,
766 20142492. (doi:10.1098/rspb.2014.2492)
- 767 5. Cleasby, I. R., Nakagawa, S. & Schielzeth, H. 2015 Quantifying the predictability of
768 behaviour: statistical approaches for the study of between-individual variation in the
769 within-individual variance. *Methods Ecol. Evol.* **6**, 27–37. (doi:10.1111/2041-
770 210X.12281)
- 771 6. Mitchell, D. J., Fanson, B. G., Beckmann, C. & Biro, P. A. 2016 Towards powerful
772 experimental and statistical approaches to study intraindividual variability in labile traits.
773 *Open Sci.* **3**, 160352. (doi:10.1098/rsos.160352)
- 774 7. Stamps, J. A., Briffa, M. & Biro, P. A. 2012 Unpredictable animals: individual
775 differences in intraindividual variability (IIV). *Anim. Behav.* **83**, 1325–1334.
776 (doi:10.1016/j.anbehav.2012.02.017)
- 777 8. Westneat, D. F., Wright, J. & Dingemanse, N. J. 2015 The biology hidden inside
778 residual within-individual phenotypic variation. *Biol. Rev.* **90**, 729–743.
779 (doi:10.1111/brv.12131)
- 780 9. DeWitt, T. J. 2016 Expanding the phenotypic plasticity paradigm to broader views of
781 trait space and ecological function. *Curr. Zool.* **62**, 463–473. (doi:10.1093/cz/zow085)

- 782 10. Scheiner, S. M. 2014 The genetics of phenotypic plasticity. XIII. Interactions with
783 developmental instability. *Ecol. Evol.* **4**, 1347–1360. (doi:10.1002/ece3.1039)
- 784 11. Tonsor, S. J., Elnaccash, T. W. & Scheiner, S. M. 2013 Developmental instability is
785 genetically correlated with phenotypic plasticity, constraining heritability, and fitness.
786 *Evolution* **67**, 2923–2935. (doi:10.1111/evo.12175)
- 787 12. Oates, A. C. 2011 What’s all the noise about developmental stochasticity?
788 *Development* **138**, 601–607. (doi:10.1242/dev.059923)
- 789 13. Fratkin, J. L., Sinn, D. L., Patall, E. A. & Gosling, S. D. 2013 Personality consistency
790 in dogs: a meta-analysis. *PLOS ONE* **8**, e54907. (doi:10.1371/journal.pone.0054907)
- 791 14. Wilsson, E. & Sundgren, P.-E. 1998 Behaviour test for eight-week old puppies -
792 heritabilities of tested behaviour traits and its correspondence to later behaviour. *Appl.*
793 *Anim. Behav. Sci.* **58**, 151–162. (doi:10.1016/S0168-1591(97)00093-2)
- 794 15. Sinn, D. L., Gosling, S. D. & Hilliard, S. 2010 Personality and performance in
795 military working dogs: reliability and predictive validity of behavioral tests. *Appl. Anim.*
796 *Behav. Sci.* **127**, 51–65. (doi:10.1016/j.applanim.2010.08.007)
- 797 16. Riemer, S., Müller, C., Virányi, Z., Huber, L. & Range, F. 2014 The predictive value
798 of early behavioural assessments in pet dogs a longitudinal study from neonates to adults.
799 *PLOS ONE* **9**, e101237. (doi:10.1371/journal.pone.0101237)
- 800 17. Kis, A., Klausz, B., Persa, E., Miklósi, Á. & Gácsi, M. 2014 Timing and presence of
801 an attachment person affect sensitivity of aggression tests in shelter dogs. *Vet. Rec.* **174**,
802 196. (doi: 10.1136/vr.101955)
- 803 18. Serpell, J. A. & Duffy, D. L. 2016 Aspects of juvenile and adolescent environment
804 predict aggression and fear in 12-month-old guide dogs. *Font. Vet. Sci.* **3**.
805 (doi:10.3389/fvets.2016.00049)
- 806 19. Robinson, L. M., Thompson, R. S. & Ha, J. C. 2016 Puppy temperament assessments
807 predict breed and American Kennel Club group but not adult temperament. *J. Appl.*
808 *Anim. Welf. Sci.* **19**, 101–114. (doi:10.1080/10888705.2015.1127765)

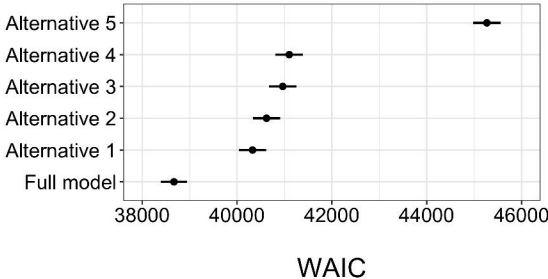
- 809 20. Marder, A. R., Shabelansky, A., Patronek, G. J., Dowling-Guyer, S. & D'Arpino, S.
810 S. 2013 Food-related aggression in shelter dogs: a comparison of behavior identified by a
811 behavior evaluation in the shelter and owner reports after adoption. *Appl. Anim. Behav.*
812 *Sci.* **148**, 150–156. (doi:10.1016/j.applanim.2013.07.007)
- 813 21. Mohan-Gibbons, H., Weiss, E. & Slater, M. 2012 Preliminary investigation of food
814 guarding behavior in shelter dogs in the United States. *Animals* **2**, 331–346.
815 (doi:10.3390/ani2030331)
- 816 22. Mornement, K. M., Coleman, G. J., Toukhsati, S. R. & Bennett, P. C. 2015
817 Evaluation of the predictive validity of the Behavioural Assessment for Re-homing K9's
818 (B.A.R.K.) protocol and owner satisfaction with adopted dogs. *Appl. Anim. Behav. Sci.*
819 **167**, 35–42. (doi:10.1016/j.applanim.2015.03.013)
- 820 23. Poulsen, A. H., Lisle, A. T. & Phillips, C. J. C. 2010 An evaluation of a behaviour
821 assessment to determine the suitability of shelter dogs for rehoming. *Vet. Med. Int.* **2010**,
822 e523781. (doi:10.4061/2010/523781)
- 823 24. Rayment, D. J., Groef, B. D., Peters, R. A. & Marston, L. C. 2015 Applied
824 personality assessment in domestic dogs: limitations and caveats. *Appl. Anim. Behav. Sci.*
825 **163**, 1–18. (doi:10.1016/j.applanim.2014.11.020)
- 826 25. Hennessy, M. B. 2013 Using hypothalamic-pituitary-adrenal measures for assessing
827 and reducing the stress of dogs in shelters: a review. *Appl. Anim. Behav. Sci.* **149**, 1–12.
828 (doi:10.1016/j.applanim.2013.09.004)
- 829 26. Protopopova, A. 2016 Effects of sheltering on physiology, immune function,
830 behavior, and the welfare of dogs. *Physiol. Behav.* **159**, 95–103.
831 (doi:10.1016/j.physbeh.2016.03.020)
- 832 27. Stephen, J. M. & Ledger, R. A. 2005 An audit of behavioral indicators of poor
833 welfare in kennelled dogs in the United Kingdom. *J. Appl. Anim. Welf. Sci.* **8**, 79–95.
834 (doi:10.1207/s15327604jaws0802_1)

- 835 28. Rooney, N. J., Gaines, S. A. & Bradshaw, J. W. S. 2007 Behavioural and
836 glucocorticoid responses of dogs (*Canis familiaris*) to kennelling: investigating
837 mitigation of stress by prior habituation. *Physiol. Behav.* **92**, 847–854.
838 (doi:10.1016/j.physbeh.2007.06.011)
- 839 29. Patronek, G. J. & Bradley, J. 2016 No better than flipping a coin: reconsidering
840 canine behavior evaluations in animal shelters. *J. Vet. Behav.* **15**, 66–77.
841 (doi:10.1016/j.jveb.2016.08.001)
- 842 30. Protopopova, A. & Wynne, C. D. L. 2014 Adopter-dog interactions at the shelter:
843 behavioral and contextual predictors of adoption. *Appl. Anim. Behav. Sci.* **157**, 109–116.
844 (doi:10.1016/j.applanim.2014.04.007)
- 845 31. Martin, J. G. A., Pirottay, E., Petellez, M. B. & Blumstein, D. T. 2017 Genetic basis
846 of between-individual and within-individual variance of docility. *J. Evol. Biol.*
847 (doi:10.1111/jeb.13048)
- 848 32. Kruschke, J. 2014 *Doing bayesian data analysis: A tutorial with r, jags, and stan.*
849 Academic Press.
- 850 33. McElreath, R. 2015 *Statistical Rethinking: A Bayesian Course with Examples in R*
851 *and Stan.* CRC Press.
- 852 34. Voith, V. L. et al. 2013 Comparison of visual and DNA breed identification of dogs
853 and inter-observer reliability. *American Journal of Sociological Research* **3**, 17–29. (doi:
854 10.1080/10888700902956151)
- 855 35. Goold, C. & Newberry, R. C. 2017 Aggressiveness as a latent personality trait of
856 domestic dogs: testing local independence and measurement invariance. *bioRxiv*
857 (doi:10.1101/117440)
- 858 36. Owczarczak-Garstecka, S. C. & Burman, O. H. 2016 Can sleep and resting
859 behaviours be used as indicators of welfare in shelter dogs (*Canis lupus familiaris*)?
860 *PLOS ONE* **11**, e0163620. (doi: <https://doi.org/10.1371/journal.pone.0163620>)

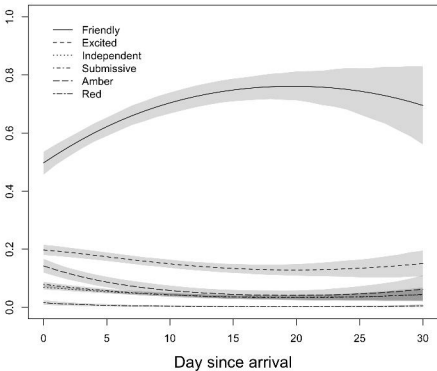
- 861 37. R Development Core Team 2016 R: A language and environment for statistical
862 computing. Vienna, Austria.
- 863 38. Tastle, W. J. & Wierman, M. J. 2007 Consensus and dissention: a measure of ordinal
864 dispersion. *Int. J. Approx. Reason.* **45**, 531–545. (doi:10.1016/j.ijar.2006.06.024)
- 865 39. Ruedin, D. 2016 agrmt: Calculate Agreement or Consensus in Ordered Rating Scales.
866 R package version 1.40.4.
- 867 40. Liddell, T. M. & Kruschke, J. K. 2015 Analyzing ordinal data: support for a Bayesian
868 approach. *SSRN*. (doi: <http://dx.doi.org/10.2139/ssrn.2692323>)
- 869 41. Foulley, J.-L. & Jaffrézic, F. 2010 Modelling and estimating heterogeneous variances
870 in threshold models for ordinal discrete data via Winbugs/Openbugs. *Comput. Methods*
871 *Programs Biomed.* **97**, 19–27. (doi:10.1016/j.cmpb.2009.05.004)
- 872 42. Kizilkaya, K. & Tempelman, R. J. 2005 A general approach to mixed effects
873 modeling of residual variances in generalized linear mixed models. *Genet. Sel. Evol.* **37**,
874 31. (doi:10.1186/1297-9686-37-1-31)
- 875 43. Nakagawa, S. & Schielzeth, H. 2010 Repeatability for Gaussian and non-Gaussian
876 data: a practical guide for biologists. *Biol. Rev.* **85**, 935–956. (doi:10.1111/j.1469-
877 185X.2010.00141.x)
- 878 44. Brommer, J. E. 2013 Variation in plasticity of personality traits implies that the
879 ranking of personality measures changes between environmental contexts: calculating the
880 cross-environmental correlation. *Behav. Ecol. Sociobiol.* **67**, 1709–1718.
- 881 45. Stan Development Team 2016 Stan modeling language users guide and reference
882 manual. Version 2.15.0.
- 883 46. Lewandowski, D., Kurowicka, D. & Joe, H. 2009 Generating random correlation
884 matrices based on vines and extended onion method. *J. Multivar. Anal.* **100**, 1989–2001.
885 (doi:10.1016/j.jmva.2009.04.008)
- 886 47. Watanabe, S. 2010 Asymptotic equivalence of Bayes cross validation and widely
887 applicable information criterion in singular learning theory. *J. Mach. Learn. Res.* **11**,

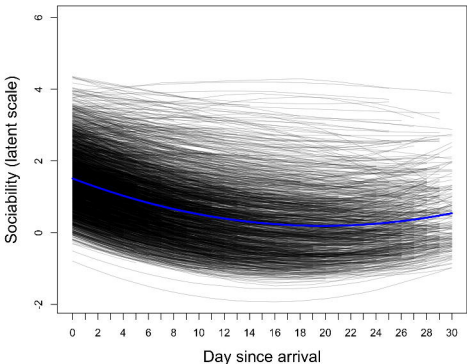
- 888 3571–3594. (url: <http://www.jmlr.org/papers/v11/watanabe10a.html>)
- 889 48. Stan Development Team 2016 Rstan: R Interface to Stan. R package version 2.14.1.
- 890 49. Hiby, E. F., Rooney, N. J. & Bradshaw, J. W. S. 2006 Behavioural and physiological
891 responses of dogs entering re-homing kennels. *Physiol. Behav.* **89**, 385–391.
892 (doi:10.1016/j.physbeh.2006.07.012)
- 893 50. Svartberg, K. & Forkman, B. 2002 Personality traits in the domestic dog (*Canis*
894 *familiaris*). *Appl. Anim. Behav. Sci.* **79**, 133–155. (doi:10.1016/S0168-1591(02)00121-1)
- 895 51. Hsu, Y. & Serpell, J. A. 2003 Development and validation of a questionnaire for
896 measuring behavior and temperament traits in pet dogs. *J. Am. Vet. Med. Assoc.* **223**,
897 1293–1300. (doi:10.2460/javma.2003.223.1293)
- 898 52. Jones, A. C. & Gosling, S. D. 2005 Temperament and personality in dogs (*Canis*
899 *familiaris*): a review and evaluation of past research. *Appl. Anim. Behav. Sci.* **95**, 1–53.
900 (doi:10.1016/j.applanim.2005.04.008)
- 901 53. Riemer, S., Müller, C., Virányi, Z., Huber, L. & Range, F. 2016 Individual and group
902 level trajectories of behavioural development in Border collies. *Appl. Anim. Behav. Sci.*
903 **180**, 78–86. (doi:10.1016/j.applanim.2016.04.021)
- 904 54. Cramer, A. O. J., Borkulo, C. D. van, Giltay, E. J., Maas, H. L. J. van der, Kendler, K.
905 S., Scheffer, M. & Borsboom, D. 2016 Major depression as a complex dynamic system.
906 *PLOS ONE* **11**, e0167490. (doi:10.1371/journal.pone.0167490)
- 907 55. Wichers, M., Groot, P. C. & Psychosystems, ESM Group, EWS Group 2016 Critical
908 slowing down as a personalized early warning signal for depression. *Psychother.*
909 *Psychosom.* **85**, 114–116. (doi:10.1159/000441458)
- 910 56. David, J. T., Cervantes, M. C., Trosky, K. A., Salinas, J. A. & Delville, Y. 2004 A
911 neural network underlying individual differences in emotion and aggression in male
912 golden hamsters. *Neuroscience* **126**, 567–578. (doi:10.1016/j.neuroscience.2004.04.031)

- 913 57. Betini, G. S. & Norris, D. R. 2012 The relationship between personality and plasticity
914 in tree swallow aggression and the consequences for reproductive success. *Anim. Behav.*
915 **83**, 137–143. (doi:10.1016/j.anbehav.2011.10.018)
- 916 58. Dewitt, T. J., Sih, A. & Wilson, D. S. 1998 Costs and limits of phenotypic plasticity.
917 *Trends Ecol. Evol.* **13**, 77–81.
- 918 59. Highcock, L. & Carter, A. J. 2014 Intraindividual variability of boldness is repeatable
919 across contexts in a wild lizard. *PLOS ONE* **9**, e95179.
920 (doi:10.1371/journal.pone.0095179)
- 921 60. Araya-Ajoy, Y. G. & Dingemanse, N. J. 2017 Repeatability, heritability, and age-
922 dependence of seasonal plasticity in aggressiveness in a wild passerine bird. *J. Anim.*
923 *Ecol.* **86**, 227–238. (doi:10.1111/1365-2656.12621)
- 924 61. Dykiert, D., Der, G., Starr, J. M. & Deary, I. J. 2012 Age differences in intra-
925 individual variability in simple and choice reaction time: systematic review and meta-
926 analysis. *PLOS ONE* **7**, e45759. (doi:10.1371/journal.pone.0045759)



Probability of sociability code





Individuals

3000

2000

1000

0

-3

-2

-1

0

1

2

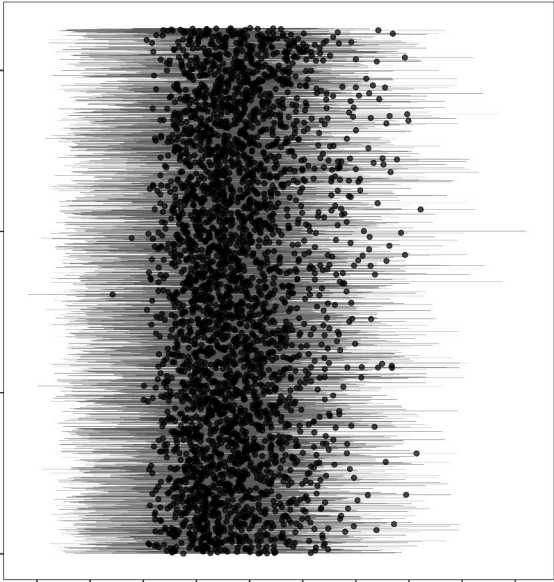
3

4

5

6

Intercepts



Individuals

3000

2000

1000

0

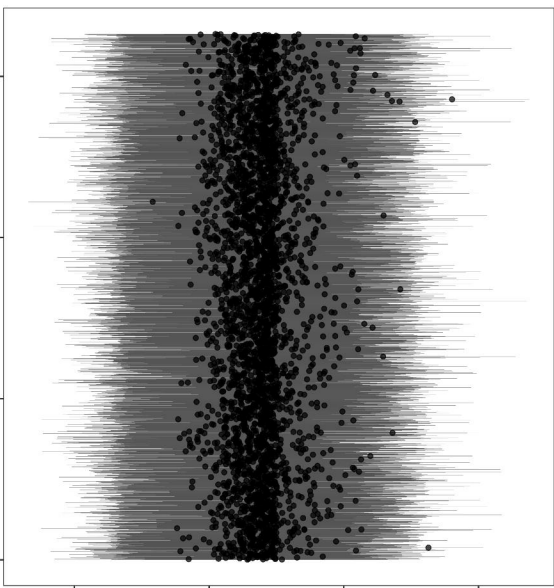
-2

-1

0

1

Linear slopes

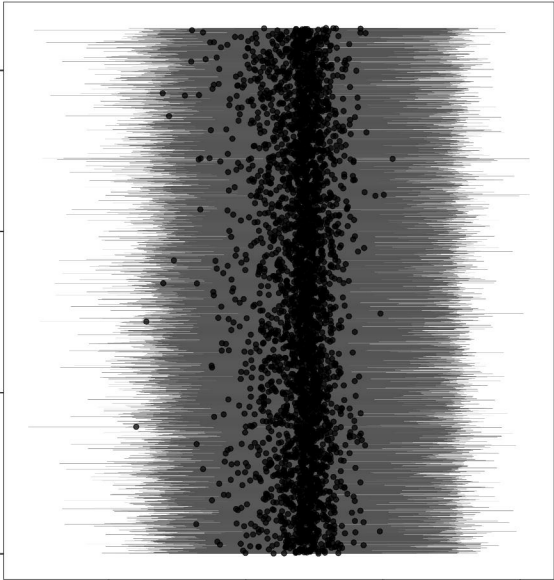


Individuals

3000
2000
1000
0

-0.5 0.0 0.5 1.0

Quadratic slopes



Individuals

3000

2000

1000

0

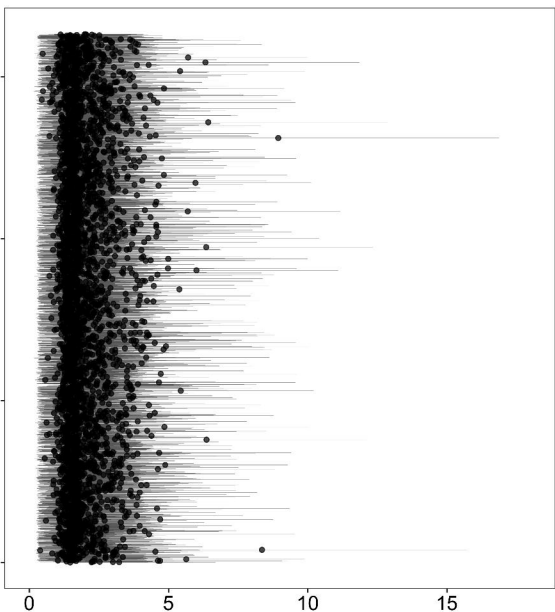
0

5

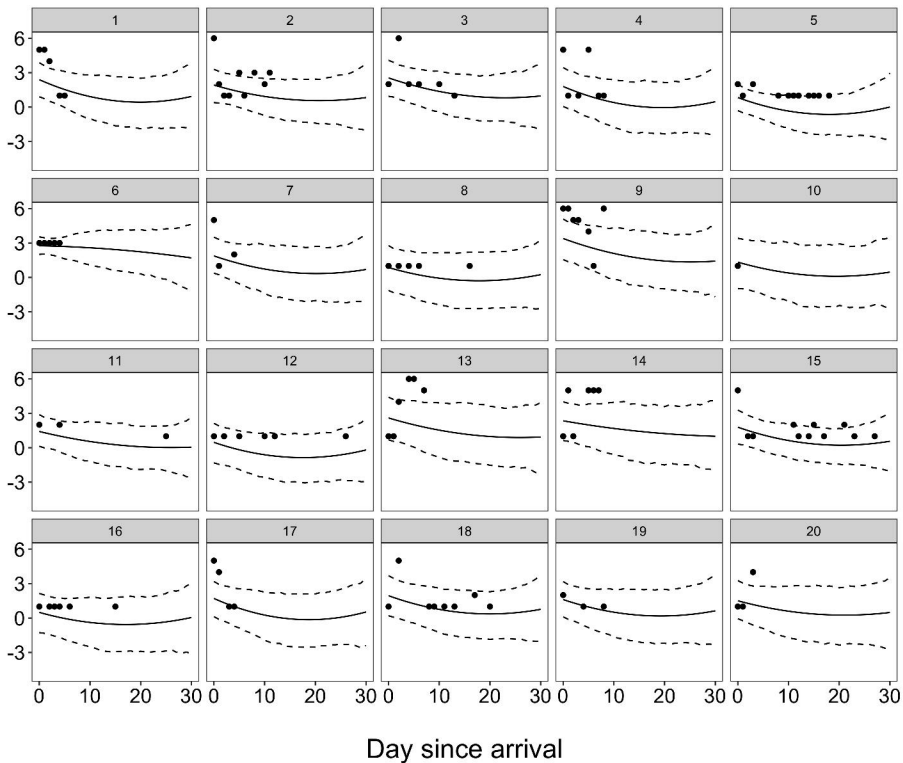
10

15

Residual SDs



Sociability (ordinal and latent scale)



Sociability (ordinal/latent scale)

