



Cite this article: Goold C, Newberry RC. 2017
Modelling personality, plasticity and
predictability in shelter dogs. *R. Soc. open sci.*
4: 170618.
<http://dx.doi.org/10.1098/rsos.170618>

Received: 5 June 2017

Accepted: 12 August 2017

Subject Category:

Biology (whole organism)

Subject Areas:

behaviour

Keywords:

inter- and intra-individual differences,
behavioural reaction norms, behavioural
repeatability, longitudinal behavioural
assessment, human–animal interactions

Author for correspondence:

Conor Goold

e-mail: conor.goold@nmbu.no

Modelling personality, plasticity and predictability in shelter dogs

Conor Goold and Ruth C. Newberry

Department of Animal and Aquacultural Sciences, Faculty of Biosciences,
Norwegian University of Life Sciences, Norway

CG, 0000-0002-9198-0889; RCN, 0000-0002-5238-6959

Behavioural assessments of shelter dogs (*Canis lupus familiaris*) typically comprise standardized test batteries conducted at one time point, but test batteries have shown inconsistent predictive validity. Longitudinal behavioural assessments offer an alternative. We modelled longitudinal observational data on shelter dog behaviour using the framework of behavioural reaction norms, partitioning variance into personality (i.e. inter-individual differences in behaviour), plasticity (i.e. inter-individual differences in average behaviour) and predictability (i.e. individual differences in residual intra-individual variation). We analysed data on interactions of 3263 dogs ($n = 19\,281$) with unfamiliar people during their first month after arrival at the shelter. Accounting for personality, plasticity (linear and quadratic trends) and predictability improved the predictive accuracy of the analyses compared to models quantifying personality and/or plasticity only. While dogs were, on average, highly sociable with unfamiliar people and sociability increased over days since arrival, group averages were unrepresentative of all dogs and predictions made at the individual level entailed considerable uncertainty. Effects of demographic variables (e.g. age) on personality, plasticity and predictability were observed. Behavioural repeatability was higher one week after arrival compared to arrival day. Our results highlight the value of longitudinal assessments on shelter dogs and identify measures that could improve the predictive validity of behavioural assessments in shelters.

1. Introduction

Personality, defined by inter-individual differences in average behaviour, represents just one component of behavioural variation of interest in animal behaviour research. Personality frequently describes less than 50% of behavioural variation in animal personality studies [1,2], leading to the combined analysis of personality with *plasticity*, individual differences in behavioural change [3], and *predictability*, individual differences in residual intra-individual variability [4–8]. These different sources of

behavioural variation can be understood using the general framework of behavioural reaction norms [3,5] that provides insight into how animals react to fluctuating environments through time and across contexts. The concept of behavioural reaction norms is built upon the use of hierarchical statistical models to quantify between- and within-individual variation in behaviour, following methods in quantitative genetics [3]. More generally, these developments reflect increasing interest across biology in expanding the 'trait space' of phenotypic evolution [9] beyond mean trait differences and systematic plasticity across environmental gradients to include residual trait variation (e.g. developmental instability [10,11]; stochastic variation in gene expression [12]).

Modest repeatability of behaviour has been documented in domestic dogs (*Canis lupus familiaris*), providing evidence for personality variation. For instance, using meta-analysis, Fratkin *et al.* [13] found an average Pearson's correlation of behaviour through time of 0.43, explaining 19% of the behavioural variance between successive time points (where the average time interval between measurements was 21 weeks). However, the goal of personality assessments in dogs is often to predict an individual dog's future behaviour (e.g. working dogs [14,15]; pet dogs [16]) and, thus, it is important not to confuse the stability of an individual's behaviour relative to the behaviour of others with stability of intra-individual behaviour. That is, individuals could vary their behaviour in meaningful ways in response to internal (e.g. ontogeny) and external (e.g. environmental) factors while maintaining differences from other individuals. When time-related change in dog behaviour has been taken into account, behavioural change at the group level has been of primary focus (e.g. [16–18]) and no studies have explored the heterogeneity of residual variance within each dog. The predominant focus on inter-individual differences and group-level patterns of behavioural change risks obscuring important individual-level heterogeneity and may partly explain why a number of dog personality assessment tools have been unreliable in predicting future behaviour [14–16,19].

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

In this paper, we demonstrate how the framework of behavioural reaction norms can be used to quantify inter- and intra-individual differences in shelter dog behaviour. To do so, we employ data on interactions of dogs with unfamiliar people from a longitudinal and observational shelter assessment. As a core feature of personality assessments, how shelter dogs interact with unknown people is of great importance. At one extreme, if dogs bite or attempt to bite unfamiliar people, they are at risk of euthanasia [29]. At the other extreme, even subtle differences in how dogs interact with potential adopters can influence adoption success [30]. Importantly, neither may all dogs react to unfamiliar people in the same way through time at the shelter nor may all dogs show the same day-to-day fluctuation of behaviour around their average behavioural trajectories. These considerations can be explored by examining behavioural reaction norms.

The analysis of behavioural reaction norms is dependent on the use of hierarchical statistical models for partitioning variance among individuals [3,5,6]. Given that ordinal data are common in behavioural research, here we illustrate how similar hierarchical models can be applied to ordinal data using a Bayesian framework (see also [31]). Apart from distinguishing inter- from intra-individual variation, we place particular emphasis on two desirable properties of the hierarchical modelling approach taken here. First, the property of *hierarchical shrinkage* [32] offers an efficacious way of making inferences about individual-level behaviour when data are highly unbalanced and potentially unrepresentative of a dog's typical behaviour. When data are sparse for certain individuals, hierarchical shrinkage means that an

Table 1. Demographic variables of dogs in the sample analysed. Mean and standard deviation (s.d.) or the number of dogs by category (*n*) are displayed.

demographic variable	mean (s.d.)/ <i>n</i>
number of observations per dog	5.9 (3.7)
days spent at the shelter	25.8 (35.0)
age (years; all at least four 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

individual's parameter estimates (e.g. intercepts) are more similar to, or shrunken towards, the group-level estimates. Second, as any prediction of future (dog) behaviour will entail uncertainty, a Bayesian approach is attractive, because we can directly obtain a probability distribution of parameter values consistent with the data (i.e. the posterior distribution) for all parameters [32,33]. By contrast, frequentist confidence intervals (CIs) are not posterior probability distributions and, thus, their interpretation is more challenging when a goal is to understand uncertainty in parameter estimates [32].

2. Material and methods

2.1. Subjects

Behavioural data on $n = 3263$ dogs from Battersea Dogs and Cats Home's longitudinal, observational assessment model were used for analysis. The data concerned all behavioural records of dogs at the shelter during 2014 (including those arriving in 2013 or departing in 2015), filtered to include all dogs: (i) at least four months of age (to ensure all dogs were treated similarly under shelter protocols, e.g. vaccinated so eligible for walks outside and kennelled in similar areas), (ii) with at least one observation during the first 31 days since arrival at the shelter, and (iii) with complete data for demographic variables to be included in the formal analysis (table 1). Because dogs spent approximately one month at the shelter on average (table 1), we focused on this period in our analyses (arrival day 0 to day 30). We did not include breed characterization due to the unreliability of using appearance to attribute breed type to shelter dogs of uncertain heritage [34].

2.2. Shelter environment

Details of the shelter environment have been presented elsewhere [35]. Briefly, the shelter was composed of three different rehoming centres (table 1): one large inner-city centre based in London (approximate capacity: 150–200 dogs), a medium-sized suburban/rural centre based in Old Windsor (approximate capacity: 100–150 dogs), and a smaller rural centre in Brands Hatch (approximate capacity: 50 dogs). Dogs considered suitable for adoption were housed in indoor kennels (typically about 4 m × 2 m, with a shelf and bedding alcove; see also [36]). Most dogs were housed individually, and given daily access to an indoor run behind their kennel. Feeding, exercising and kennel cleaning were performed by a relatively stable group of staff members. Dogs received water *ad libitum* and two meals daily according to veterinary recommendations. Sensory variety was introduced daily (e.g. toys, essential oils, classical music, access to quiet 'chill-out' rooms). Regular work hours were from 08.00 to 17.00 each day, with public visitation from 1000 to 1600 h. Dogs were socialized with staff and/or volunteers daily.

2.3. Data collection

The observational assessment implemented at the shelter included observations of dogs by trained shelter employees in different, everyday contexts, each with its own qualitative ethogram of possible

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

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 does not move away from the person
8. stressed	amber	0.5	high frequency/intensity of stress behaviours, which may include dribbling, stereotypic behaviours, stress vocalizations, 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

behaviours. Shortly after dogs were observed in relevant contexts, employees entered observations into a custom, online platform using computers located in different housing areas. Each behaviour within a context had its own code. Previously, we have reported on aggressive behaviour across contexts [35]. Here, we focus on variation in behaviour in one of the most important contexts, 'Interactions with unfamiliar people', which pertained to how dogs reacted when people with whom they had never interacted before approached, made eye contact, spoke to and/or attempted to make physical contact with them. For the most part, this context occurred outside of the kennel, but it could also occur if an unfamiliar person entered the kennel. Observations could be recorded by an employee meeting an unfamiliar dog, or by an employee observing a dog meeting an unfamiliar person. Different employees could input records for the same dog, and employees could discuss the best code to describe a certain observation if required.

Behavioural observations in the 'Interactions with unfamiliar people' context were recorded using a 13-code ethogram (table 2). Each behavioural code was subjectively labelled and generally defined, providing a balance between behavioural rating and behavioural coding methodologies. The ethogram represented a scale of behavioural problem severity and assumed adoptability (higher codes indicating higher severity of problematic behaviour/lower sociability), reflected by grouping the 13 codes further into green, amber and red codes (table 2). Green behaviours posed no problems for adoption, amber behaviours suggested dogs may require some training to facilitate successful adoption, but did not pose a danger to people or other dogs, and red behaviours suggested dogs needed training or behavioural modification to facilitate successful adoption and could pose a risk to people or other dogs. A dog's suitability for adoption was, however, based on multiple behavioural observations over a number of

days. When registering an observation, the employee selected the highest code in the ethogram that was observed on that occasion (i.e. the most severe level of problematic behaviour was given priority). There were periods when a dog could receive no entries for the context for several days, but other times when multiple observations were recorded on the same day, usually when a previous observation was followed by a more serious behavioural event. In these instances, and in keeping with the shelter protocol, we retained the highest (i.e. most severe) behavioural code registered for the context that day. When the behaviours were the same, only one record was retained for that day. This resulted in an average of 5.9 (s.d. = 3.7; range = 1–22) records per dog on responses during interactions with unfamiliar people while at the shelter. For dogs with more than one record, the average number of days between records was 2.8 (s.d. = 2.2; range = 1–29).

2.4. Validity and inter-rater reliability

Inter-rater reliability and the validity of the assessment methodology were evaluated using data from a larger research project at the shelter. Videos depicting different behaviours in different contexts were filmed by canine behaviourists working at the shelter, who subsequently organized video coding sessions with 93 staff members (each session with about 5–10 participants) across rehoming centres [35]. The authors were blind to the videos and administration of video coding sessions. The staff members were shown 14 videos (each about 30 s long) depicting randomly selected behaviours, two from each of seven different assessment contexts (presented in a pseudo-random order, the same for all participants). Directly after watching each video, they individually recorded (on a paper response form) which ethogram code best described the behaviour observed in each context. Two videos depicted behaviour during interactions with people (familiar versus unfamiliar not differentiated), one demonstrating *Reacts to people aggressive* and the other *Reacts to people non-aggressive* (table 2). Below, we present the inter-rater reliabilities and the percentage of people who chose the correct behaviour and colour category for these two videos in particular, but also the averaged results across the 14 videos, because there was some redundancy between ethogram scales across contexts.

2.5. Statistical analyses

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

2.5.1. Validity and inter-rater reliability

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

2.5.2. Hierarchical Bayesian ordinal probit model

The distribution of ethogram categories was heavily skewed in favour of the green codes (table 2), particularly the first *Friendly* category. As some categories were chosen particularly infrequently, we aggregated the raw responses into a 6-category scale: (i) *Friendly*, (ii) *Excitable*, (iii) *Independent*, (iv) *Submissive*, (v) *Amber codes*, and (vi) *Red codes*. This aggregated scale retained the main variation in the data and simplified the data interpretation. We analysed the data using a Bayesian ordinal probit model (described in [32,40]), but extended to integrate the hierarchical structure of the data, including heteroscedastic residual standard deviations, to quantify predictability for each dog (for related models, see [31,41,42]). The ordinal probit model, also known as the cumulative or thresholded normal model, is motivated by a latent variable interpretation of the ordinal scale. That is, an ordinal dependent variable,

Y , with categories K_j , from $j = 1$ to J , is a realization of an underlying continuous variable divided into thresholds, θ_c , for $c = 1$ to $J - 1$. Under the probit model, the probability of each ordinal category is equal to its area under the cumulative normal distribution, Φ , with mean, μ , s.d. σ and thresholds θ_c :

$$\text{Prob}(Y = K | \mu, \sigma, \theta_c) = \Phi\left[\frac{\theta_c - \mu}{\sigma}\right] - \Phi\left[\frac{\theta_{c-1} - \mu}{\sigma}\right]. \quad (2.1)$$

For the first and last categories, this simplifies to $\Phi[(\theta_c - \mu)/\sigma]$ and $1 - \Phi[(\theta_{c-1} - \mu)/\sigma]$, respectively. As such, the latent scale extends from $\pm\infty$. Here, the ordinal dependent variable was a realization of the hypothesized continuum of ‘insociability when meeting unfamiliar people’, with six categories and five threshold parameters. While ordinal regression models usually fix the mean and s.d. of the latent scale to 0 and 1 and estimate the threshold parameters, we fixed the first and last thresholds to 1.5 and 5.5, respectively, allowing for the remaining thresholds, and the mean and s.d., to be estimated from the data. As explained by Kruschke [32], this allows for the results to be interpretable with respect to the ordinal scale. We present the results using both the predicted probabilities of ordinal sociability codes and estimates on the latent, unobserved scale assumed to generate the ordinal responses.

2.5.3. Hierarchical structure

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

$$y_{ij}^* = \beta_0 + \nu_{0j} + \sum_{p=1}^P \beta_{p0} x_{pj} + \left(\beta_1 + \nu_{1j} + \sum_{p=1}^P \beta_{p1} x_{pj} \right) \text{day}_{ij} + \left(\beta_2 + \nu_{2j} + \sum_{p=1}^P \beta_{p2} x_{pj} \right) \text{day}_{ij}^2. \quad (2.2)$$

The above equation expresses the longitudinal pattern of behaviour as a function of (i) a group-level intercept the same for all dogs, β_0 , and the deviation from the group-level intercept for each dog, ν_{0j} , (ii) a linear effect of day since arrival, β_1 , and each dog’s deviation, ν_{1j} , and (iii) a quadratic effect of day since arrival, β_2 , and each dog’s deviation, ν_{2j} . A quadratic effect was chosen based on preliminary plots of the data at the group level and at the individual level, although we also compared the model’s predictive accuracy with simpler models (described below). Day since arrival was standardized, meaning that the intercepts reflected the behaviour on the average day since arrival across dogs (approx. day 8). The three dog-level parameters, ν_j , correspond to personality and linear and quadratic plasticity parameters. The terms $\sum_{p=1}^P \beta_p x_{pj}$ denote the effect of P dog-level predictor variables (x_p), included to explain variance between dog-level intercepts and slopes. These included: the number of observations for each dog, the number of days dogs spent at the shelter controlling for the number of observations (i.e. the residuals from a linear regression of total number of days spent at the shelter on the number of observations), average age while at the shelter, average weight at the shelter, sex, neuter status, source type and rehoming centre (table 1). For neuter status, we did not make comparisons between the ‘undetermined’ category and other categories. The primary goal of including these predictor variables was to obtain estimates of individual differences conditional on relevant inter-individual differences variables, because the data were observational.

The s.d. model was

$$\sigma = \exp\left(\delta + \nu_{3j} + \sum_{p=1}^P \beta_{p3} x_{pj}\right). \quad (2.3)$$

This equation models the s.d. of the latent scale by its own regression, with group-level s.d. intercept, δ , evaluated at the average day, the deviation for each dog from the group-level s.d. intercept, ν_{3j} , and predictor variables, $\sum_{p=1}^P \beta_{p3} x_{pj}$, as in the mean model (equation (2.2)). The s.d.s across dogs were assumed to approximately follow a log-normal distribution, with $\ln(\sigma)$ approximately normally distributed (hence the exponential inverse-link function). The parameter ν_{3j} corresponds to each dog’s residual s.d. or predictability.

All four dog-level parameters were assumed to be multivariate normally distributed with means 0 and variance–covariance matrix Σ_v estimated from the data:

$$\Sigma_v = \begin{bmatrix} \tau_{v_0}^2 & \rho_{v_{01}} \tau_{v_0} \tau_{v_1} & \rho_{v_{02}} \tau_{v_0} \tau_{v_2} & \rho_{v_{03}} \tau_{v_0} \tau_{v_3} \\ \cdots & \tau_{v_1}^2 & \rho_{v_{12}} \tau_{v_1} \tau_{v_2} & \rho_{v_{13}} \tau_{v_1} \tau_{v_3} \\ \cdots & \cdots & \tau_{v_2}^2 & \rho_{v_{23}} \tau_{v_2} \tau_{v_3} \\ \cdots & \cdots & \cdots & \tau_{v_3}^2 \end{bmatrix}. \quad (2.4)$$

The diagonal elements are the variances of the dog-level intercepts, linear slopes, quadratic slopes and residual s.d.s, while the covariances fill the off-diagonal elements (only the upper triangle shown), where ρ is the correlation coefficient. In the results, we report τ_{v_3} (the s.d. of dog-level residual s.d.s) on the original scale, rather than the log-transformed scale, using $\sqrt{e^{2\delta + \tau_{v_3}^2} e^{\tau_{v_3}^2} - 1}$. Likewise, δ was transformed to the median of the original scale by e^δ .

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

$$\text{ICC}_i = \frac{\tau_{v_0}^2 + 2\text{Cov}_{v_0, v_1} \text{Day}_i + \tau_{v_1}^2 \text{Day}_i^2 + 2\text{Cov}_{v_0, v_2} \text{Day}_i^2 + \tau_{v_2}^2 \text{Day}_i^4 + 2\text{Cov}_{v_1, v_2} \text{Day}_i^3}{\text{numerator} + e^\delta}. \quad (2.5)$$

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

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

$$\Omega = \Psi \mathbf{K} \Psi'. \quad (2.6)$$

In equation (2.6), \mathbf{K} is the variance–covariance matrix of the dog-level intercepts and (linear and quadratic) slopes, and Ψ is a three-by-three matrix with a column vector of 1’s, a column vector containing -1 , 0 and 1 defining the day values for the cross-environmental correlations for the linear component, and a column vector containing 1 , 0 and 1 defining the day values for the cross-environmental correlations for the quadratic component. Once defined, Ω was scaled to a correlation matrix. Finally, to summarize the degree of individual differences in predictability, we calculated the ‘coefficient of variation for predictability’ as $\sqrt{e^{\tau_{v_3}^2} - 1}$ following Cleasby *et al.* [5].

2.5.4. Prior distributions

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

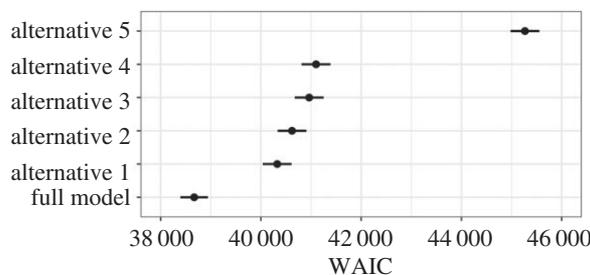


Figure 1. Out-of-sample predictive accuracy (lower is better) for each model (described in §2.5.5) measured by the WAIC. Black points denote the WAIC estimate and horizontal lines show WAIC estimates \pm s.e. Mean \pm s.e.: full model = 38 669 \pm 275; alternative 1 = 40 326 \pm 288; alternative 2 = 40 621 \pm 288; alternative 3 = 40 963 \pm 289; alternative 4 = 41 100 \pm 289; alternative 5 = 45 268 \pm 289.

2.5.5. Model selection and computation

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

Models were computed using the probabilistic programming language Stan [45] using the *RStan* package [48] v. 2.15.1, which employs Markov chain Monte Carlo estimation using Hamiltonian Monte Carlo (see the R script file and Stan code for full details). We ran four chains of 5000 iterations each, discarding the first 2500 iterations of each chain as warm-up, and setting thinning to 1. Convergence was assessed visually using trace plots to ensure chains were well mixed, numerically using the Gelman–Rubin statistic (values close to 1 and less than 1.05 indicating convergence) and by inspecting the effective sample size of each parameter. We also used graphical posterior predictive checks to assess model predictions against the raw data, including ‘counterfactual’ predictions [33] to inspect how dogs would be predicted to behave across the first month of being in the shelter regardless of their actual number of observations or length of stay at the shelter. To summarize parameter values, we calculated mean (denoted β) and 95% highest density intervals (HDIs), the 95% most probable values for each parameter (using functions in the *rethinking* package [33]). For comparing levels of categorical variables, the 95% HDIs of their differences were calculated (i.e. the differences between the coefficients at each step in the Markov chain Monte Carlo chain, denoted β_{diff}). When the 95% HDIs of predictor variables surpassed zero, a credible effect was inferred.

3. Results

3.1. Inter-rater reliability and validity

For the two videos depicting interactions with people, consensus was 0.75 (95% CI: 0.66, 0.84) for the video showing an example of *Reacts to people non-aggressive* and 0.77 (95% CI: 0.74, 0.81) for the example of *Reacts to people aggressive*. Neither did these results overlap with the null distributions (see the electronic supplementary material, table S1), indicating significant inter-rater reliability. For the video showing *Reacts to people non-aggressive*, 77% chose the correct code and 83% a code of the correct colour category (amber), and, as previously reported by Goold & Newberry [35], 52% chose the correct code for the video showing *Reacts to people aggressive* and 55% chose a code of the correct colour category (red; 42% chose

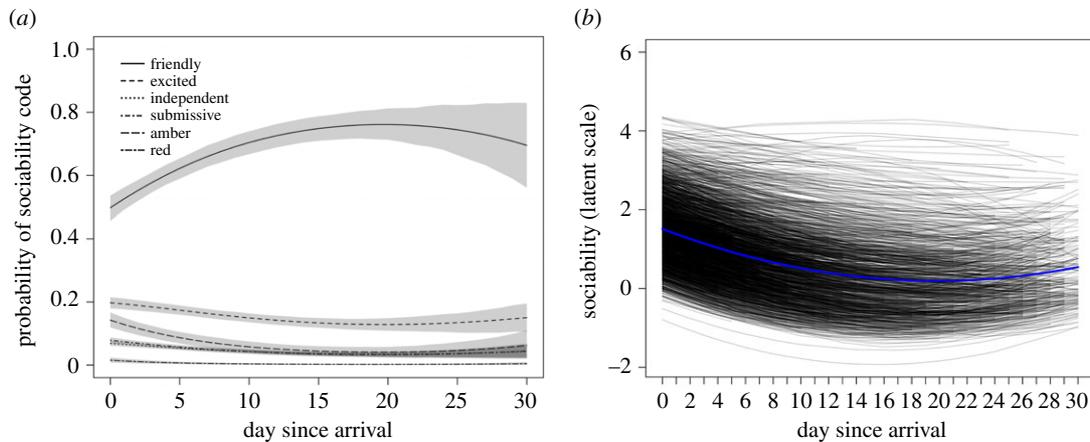


Figure 2. (a) Predicted probabilities (posterior means = black lines; 95% HDIs = shaded areas) of different sociability codes across days since arrival. (b) Posterior mean behavioural trajectories on the latent scale (ranging from $\pm\infty$) at the group level (blue line) and for each individual (black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

the amber code *Reacts to people non-aggressive* instead). Across all assessment context videos, the average consensus was 0.71 and participants chose the correct ethogram category 66% of the time, while 78% of answers were a category of the correct ethogram colour.

3.2. Hierarchical ordinal probit model

The full model had the best out-of-sample predictive accuracy, with the inclusion of heterogeneous residual s.d.s among dogs improving model fit by over 1500 WAIC points compared to the second most plausible model (alternative 1 in figure 1). In general, models that included more parameters to describe personality, plasticity and predictability, and models with a quadratic effect of day, had better out-of-sample predictive accuracy, despite the added complexity brought by additional parameters.

At the group level, the *Friendly* code (table 2) was most probable overall and was estimated to increase in probability across days since arrival, while the remaining sociability codes either decreased or stayed at low probabilities (figure 2a), reflecting the raw data. On the latent sociability scale (figure 2b), the group-level intercept parameter on the average day was 0.68 (95% HDI: 0.51, 0.86). A 1 s.d. increase in the number of days since arrival was associated with a -0.63 unit (95% HDI: -0.77 , -0.50) change on the latent scale on average (i.e. reflecting increasing sociability), and the group-level quadratic slope was positive ($\beta = 0.20$, 95% HDI: 0.10, 0.30), reflecting a quicker rate of change in sociability earlier after arrival to the shelter than later (i.e. a concave down parabola). There was a slight increase in the quadratic curve towards the end of the one-month period, although there were fewer behavioural observations at this point and so greater uncertainty about the exact shape of the curve, resulting in estimates being pulled closer to those of the intercepts. The group-level residual standard deviation had a median of 1.84 (95% HDI: 1.67, 2.02).

At the individual level, heterogeneity existed in behavioural trajectories across days since arrival (figure 2b). The s.d.s of dog-varying parameters were: (i) intercepts: 1.29 (95% HDI: 1.18, 1.41; figure 3a), (ii) linear slopes: 0.56 (95% HDI: 0.47, 0.65; figure 3b), (iii) quadratic slopes: 0.28 (95% HDI: 0.20, 0.35; figure 3c), and (iv) residual s.d.s: 1.39 (95% HDI: 1.22, 1.58; figure 3d). There was also large uncertainty in individual-level estimates. Figure 4 displays counterfactual model predictions for 20 randomly sampled dogs. Uncertainty in reaction norm estimates, illustrated by the width of the 95% HDIs (dashed black lines), was greatest when data were sparse (e.g. towards the end of the one-month study period). Hierarchical shrinkage meant that individuals with observations of less sociable responses, or individuals with few behavioural observations, tended to have model predictions pulled towards the overall mean. Note that regression lines depict values on the latent scale predicted to generate observations on the ordinal scale, and so may not clearly fit the ordinal data points. The coefficient of variation for predictability was 0.64 (95% HDI: 0.58, 0.70). Individuals with the five highest and lowest residual s.d. estimates are shown in figure 5.

Dog-varying intercepts positively correlated with linear slope parameters ($\rho = 0.38$, 95% HDI: 0.24, 0.50) and negatively correlated with quadratic slope parameters ($\rho = -0.54$, 95% HDI: -0.68 , -0.39), and

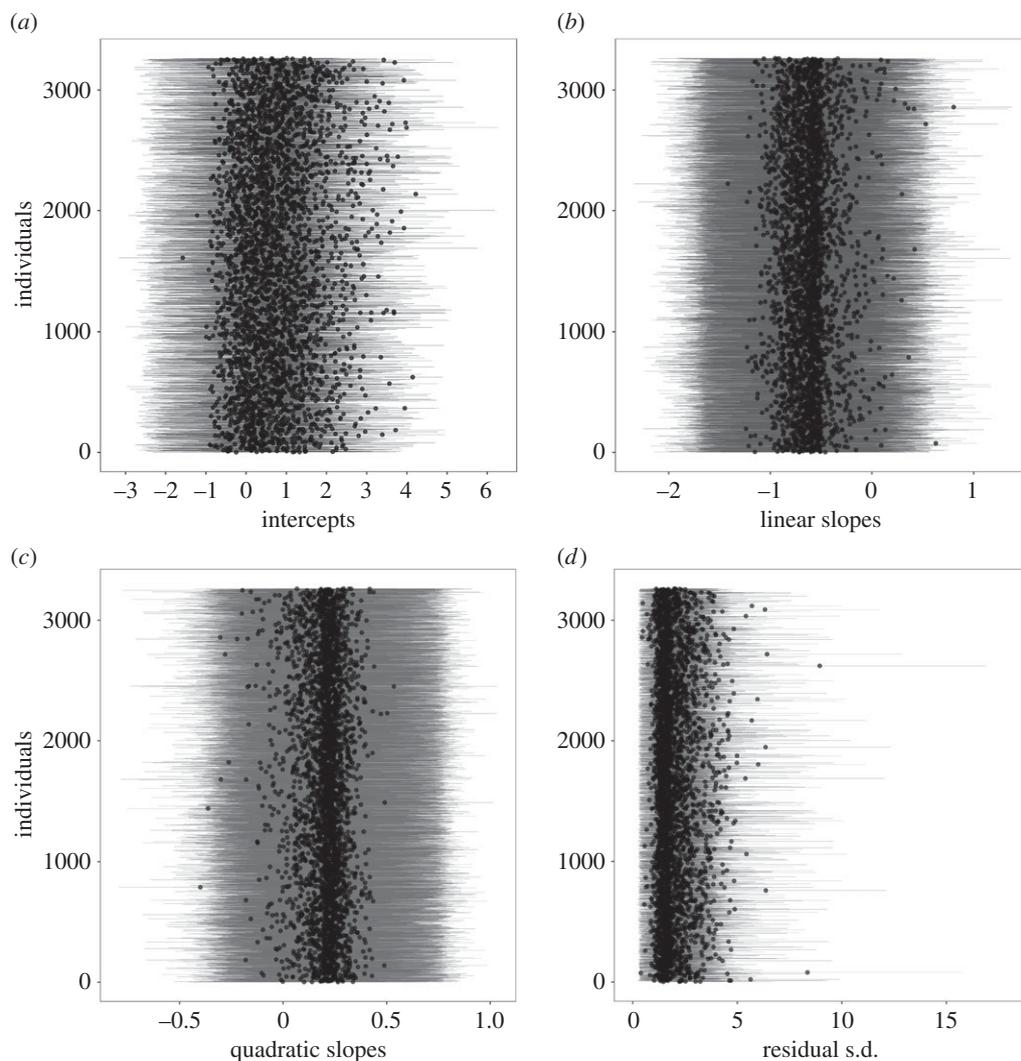


Figure 3. Posterior means (black dots) and 95% HDIs (grey horizontal bars) for each dog's (a) intercept, (b) linear slope, (c) quadratic slope and (d) residual s.d. parameter.

linear and quadratic slopes had a negative correlation ($\rho = -0.75$, 95% HDI: -0.88 , -0.59), indicating that less sociable individuals (with higher scores on the ordinal scale) had flatter reaction norms on average. Dog-varying residual s.d.s had a correlation with the intercept parameters of approximately zero ($\rho = 0.00$, 95% HDI: -0.10 , 0.10) but were negatively correlated with the linear slope parameters ($\rho = -0.37$, 95% HDI: -0.51 , -0.22) and positively correlated with the quadratic slopes ($\rho = 0.24$, 95% HDI: 0.05 , 0.42), indicating that dogs with greater residual s.d.s were predicted to change the most across days since arrival.

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

A 1 s.d. increase in the number of observations was associated with higher intercepts ($\beta = 0.12$, 95% HDI: 0.03, 0.21; see the electronic supplementary material, table S2) and higher residual s.d.s ($\beta = 0.06$, 95% HDI: 0.02, 0.10). Increasing age by 1 s.d. was associated with lower intercepts ($\beta = -0.61$, 95% HDI: -0.70 , -0.51), steeper linear slopes ($\beta = -0.20$, 95% HDI: -0.27 , -0.13), a stronger quadratic curve ($\beta = 0.07$, 95% HDI: 0.03, 0.12) and larger residual s.d.s ($\beta = 0.05$, 95% HDI: 0.01, 0.09). Increasing weight by 1 s.d. was associated with shallower quadratic curves ($\beta = -0.05$, 95% HDI: -0.09 , -0.01). No credible effect of sex was observed on personality, plasticity or predictability. Gift dogs had larger intercepts than returned dogs ($\beta_{\text{diff}} = 0.28$, 95% HDI: 0.04, 0.52) and stray dogs ($\beta_{\text{diff}} = 0.33$, 95% HDI: 0.15, 0.50), as well as steeper linear slopes ($\beta_{\text{diff}} = -0.25$, 95% HDI: -0.38 , -0.13) and higher residual s.d.s than

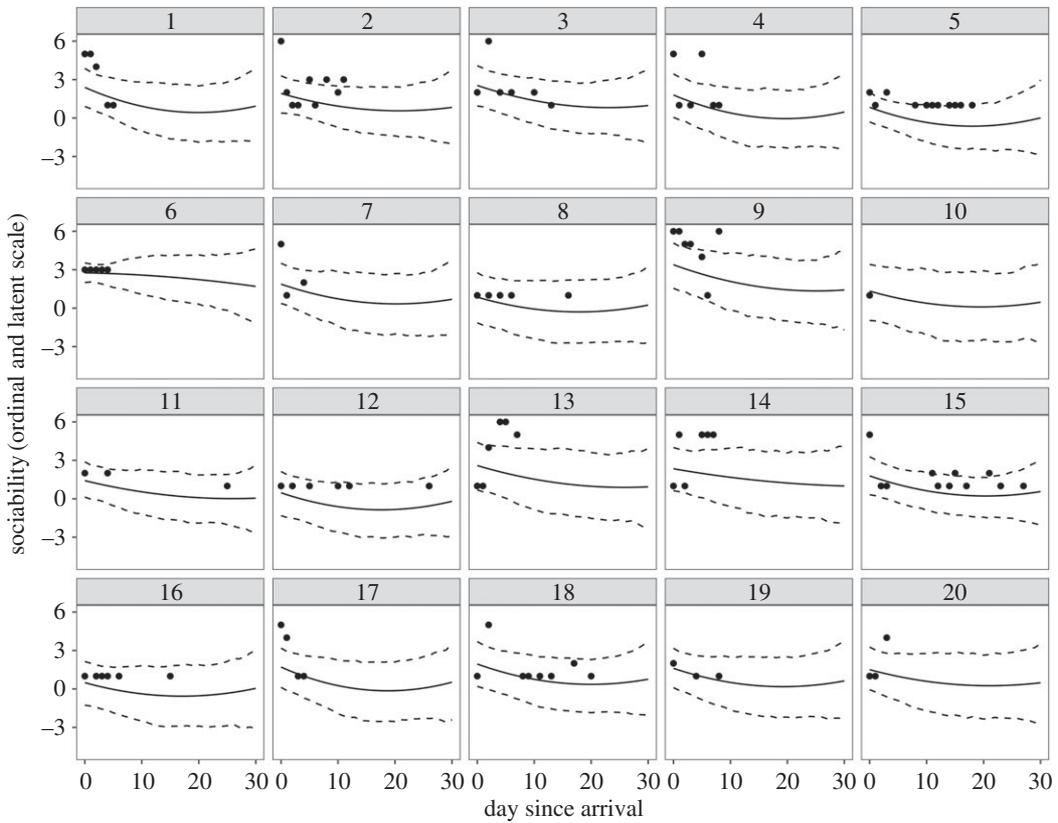


Figure 4. Predicted reaction norms ('counterfactual' plots) for 20 randomly selected dogs. Black points show raw data on the ordinal scale (higher values indicate lower sociability), and solid and dashed lines illustrate posterior means and 95% HDIs. When data were sparse, there was increased uncertainty in model predictions. Owing to the hierarchical shrinkage, model predictions of individual dogs were pulled towards the group-level mean, particularly for those dogs showing higher behavioural codes (i.e. less sociable responses).

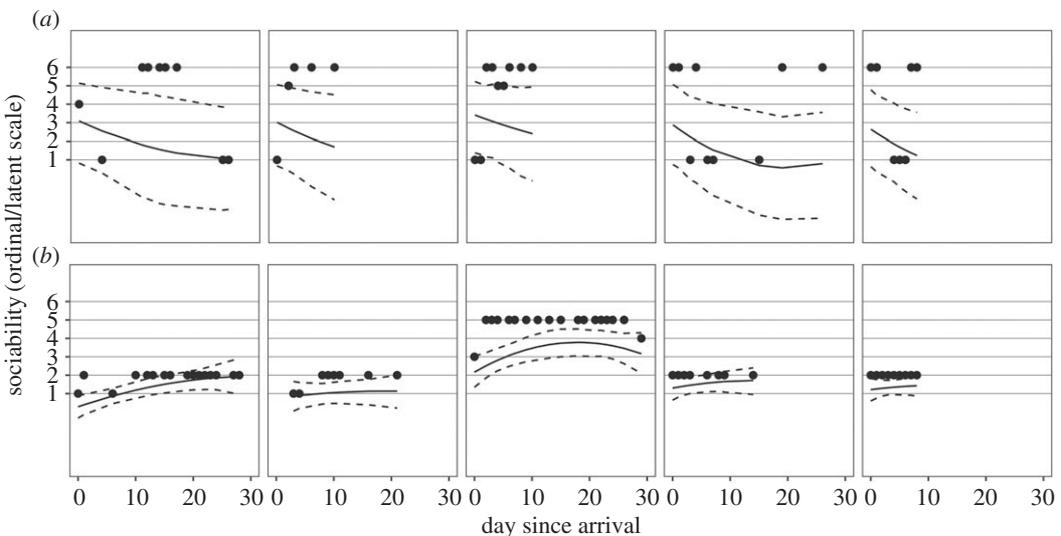


Figure 5. Reaction norms (posterior means = solid black lines; 95% HDIs = dashed black lines) for individuals with the five highest (a) and five lowest (b) residual s.d.s. Black points represent raw data on the ordinal scale (higher values indicating lower sociability).

stray dogs ($\beta_{\text{diff}} = 0.10$, 95% HDI: 0.02, 0.18). Dogs at the large rehoming centre had steeper linear slopes ($\beta_{\text{diff}} = -0.70$, 95% HDI: -0.84 , -0.56) and stronger quadratic curves ($\beta_{\text{diff}} = 0.35$, 95% HDI: 0.26, 0.45) than dogs at the medium rehoming centre, and lower intercept parameters ($\beta_{\text{diff}} = -0.30$, 95% HDI: -0.50 , -0.09) and steeper linear slopes ($\beta_{\text{diff}} = -0.22$, 95% HDI: -0.38 , -0.06) than dogs at the small

rehoming centre. Compared to dogs at the small rehoming centre, dogs at the medium centre had lower intercepts ($\beta_{\text{diff}} = -0.25$, 95% HDI: $-0.48, -0.01$), and shallower linear ($\beta_{\text{diff}} = 0.48$, 95% HDI: $0.30, 0.66$) and quadratic slopes ($\beta_{\text{diff}} = -0.34$, 95% HDI: $-0.46, -0.22$). Dogs already neutered before arrival to the shelter had lower intercepts ($\beta_{\text{diff}} = -0.54$, 95% HDI: $-1.07, -0.03$) and lower residual s.d.s ($\beta_{\text{diff}} = -0.53$, 95% HDI: $-0.85, -0.22$) than dogs not neutered, but higher intercepts ($\beta_{\text{diff}} = 0.20$, 95% HDI: $0.03, 0.37$) and higher residual s.d.s ($\beta_{\text{diff}} = 0.10$, 95% HDI: $0.02, 0.19$) than those neutered while at the shelter. Unneutered dogs had higher intercepts ($\beta_{\text{diff}} = 0.74$, 95% HDI: $0.20, 1.26$) and higher residual s.d.s ($\beta_{\text{diff}} = 0.63$, 95% HDI: $0.30, 0.92$) than dogs neutered at the shelter.

4. Discussion

This study applied the framework of behavioural reaction norms to quantify inter- and intra-individual differences in shelter dog behaviour during interactions with unfamiliar people. This is the first study to systematically analyse behavioural data from a longitudinal, observational assessment of shelter dogs. Dogs demonstrated substantial individual differences in personality, plasticity and predictability, which were not well described by simply investigating how dogs behaved on average. In particular, accounting for individual differences in predictability, or the short-term, day-to-day fluctuations in behaviour, resulted in significant improvement in model fit (figure 1). The longitudinal modelling of dog behaviour also demonstrated that behavioural repeatability increased with days since arrival (i.e. increasing proportion of variance explained by between-individual differences), particularly across the first week since arrival. Similarly, while individuals maintained rank-order differences in sociability across smaller periods (i.e. first 8 days), rank-order differences were only moderately maintained between arrival at the shelter and day 15. The results highlight the importance of adopting observational and longitudinal assessments of shelter dog behaviour, provide a method by which to analyse longitudinal data commensurate with other work in animal behaviour, and identify previously unconsidered behavioural measures that could be used to improve the predictive validity of behavioural assessments in dogs.

4.1. Average behaviour

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

Previous studies provide conflicting evidence regarding how shelter dogs adapt to the kennel environment over time, including behavioural and physiological profiles indicative of both positive and negative welfare [26]. Whereas some authors report decreases in the prevalence of some stress- and/or fear-related behaviour with time [27,49], others have reported either no change or an increase in behaviours indicative of poor welfare [17,30]. Of relevance here, Kis *et al.* [17] found that aggression towards unknown people increased over the first two weeks of being at a shelter. In the current study, aggression was rare (table 2), and the probability of ‘red codes’ (which included aggression) decreased with days at the shelter (figure 3a). A salient difference is that Kis *et al.* [17] collected data using a standardized behavioural test consisting of a stranger engaging in a ‘threatening approach’ towards dogs. By contrast, we used a large data set of behavioural observations recorded after non-standardized, spontaneous interactions between dogs and unfamiliar people. In recording spontaneous interactions, the shelter aimed to elicit behaviour more representative of a dog’s typical behaviour outside of the shelter environment than would be seen in a standardized behavioural assessment. Previously, authors have noted that standardized behavioural assessments may induce stress and inflate the chances of dogs displaying aggression [29], emphasizing the value of observational methods of assessment in shelters [24]. While such observational methods are less standardized, they may have greater ecological

validity by giving results more representative of how dogs will behave outside of the shelter. Testing the predictive value of observational assessments on behaviour post-adoption is the focus of ongoing research.

4.2. Individual-level variation

When behavioural data are aggregated across individuals, results may provide a poor representation of how individuals in a sample actually behaved. Here, we found heterogeneity in dog behaviour across days since arrival, even after taking into account a number of dog-level predictor variables that could explain inter-individual differences. Variation in average behaviour of individuals across days (i.e. variation in intercept estimates of dogs) illustrated that personality estimates spanned a range of behavioural codes, although model predictions mostly spanned the green codes (figure 2*b* and table 2). However, while there were many records to inform group-level estimates, there were considerably fewer records available for each individual, which resulted in large uncertainty of individual personality parameters (illustrated by wide 95% HDI bars in figure 3*a*). Personality variation has been the primary focus of previous analyses of individual differences in dogs, often based on data collected at one time point and usually on a large number of behavioural variables consolidated into composite or latent variables (e.g. [50–52]). Our results highlight that ranking individuals on personality dimensions from few observations entails substantial uncertainty.

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

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

Personality and plasticity were correlated, with dogs with less sociable behaviour across days being less plastic. Previous studies have explored the relationship between how individuals behave on average and their degree of behavioural change. David *et al.* [56] found that male golden hamsters (*Mesocricetus auratus*) showing high levels of aggression in a social intruder paradigm were slower in adapting to a delayed-reward paradigm. In practice, the relationship between personality and plasticity is probably context dependent. Betini & Norris [57] found, for instance, that more aggressive male tree swallows (*Tachycineta bicolor*) during nest defence were more plastic in response to variation in temperature, but that plasticity was only advantageous for non-aggressive males and no relationship was present between personality and plasticity in females. The correlation between personality and plasticity indicates a ‘fanning out’ shape of the reaction norms through time (figure 2*b*). Consequently, behavioural repeatability or the amount of variance explained by between-individual differences increased as a function of day, but only after the first week after arrival. The ‘cross-environmental’ correlation,

moreover, indicated that the most sociable dogs on arrival day were not necessarily the most sociable on later days at the shelter. In particular, the correlation between sociability scores on arrival day and day 15 was only moderate, supporting Brommer [44] that the rank-ordering of trait scores is not always reliable. By contrast, the cross-environmental correlations between day 0 and 8, and between day 8 and 15, were much stronger. These results suggest that shelters using standardized behavioural assessments would benefit from administering such tests as late as possible after dogs arrive.

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

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

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

4.3. Predictors of individual variation

Finally, we found associations between certain predictor variables and personality, plasticity and predictability (electronic supplementary material, table S2). Our primary reason for including these predictor variables was to obtain more accurate estimates of personality, plasticity and predictability, and we remain cautious about *a posteriori* interpretations of their effects, especially because the theory underlying why individuals may, for example, demonstrate differences in predictability is in its infancy [8]. The reproducibility of a number of the results would, nevertheless, be interesting to confirm in future research. In particular, understanding factors affecting intra-individual change is important given that many personality assessments are used to predict an individual's future behaviour, rather than understand inter-individual differences. Here, increasing age was associated with greater plasticity (linear and quadratic change) and lower predictability, although some of the 95% HDIs of parameters were close to zero, indicative of small effects. In great tits (*Parus major*) conversely, plasticity decreased with age [60], while in humans, intra-individual variability in reaction times increased with age [61].

Moreover, non-neutered dogs showed lower predictability than neutered dogs, and dogs entering the shelter as gifts (relinquished by their owners) had lower predictability estimates than stray dogs (dogs brought in by local authorities or members of the public after being found without their owners). These results can be used to formulate specific hypotheses about behavioural variation.

5. Conclusion

We applied the framework of behavioural reaction norms to data from a longitudinal and observational shelter dog behavioural assessment, quantifying inter- and intra-individual behavioural variation in interactions of dogs with unfamiliar people. Overall, shelter dogs were sociable with unfamiliar people and sociability continued to increase with days since arrival to the shelter. At the same time, dogs showed individual differences in personality, plasticity and predictability. Accounting for all of these components substantially improved model fit, particularly the inclusion of predictability, which suggests that individual differences in day-to-day behavioural variation represent an important, yet largely unstudied, component of dog behaviour. Our results also highlight the uncertainty of making predictions about shelter dog behaviour, particularly when the number of behavioural observations is low. For shelters conducting standardized behavioural assessments, assessments are probably best carried out as late as possible, given that rank-order differences between individuals on arrival and at day 15 were only moderately related. In conclusion, this study supports moving towards observational and longitudinal assessments of shelter dog behaviour, has demonstrated a Bayesian method by which to analyse longitudinal data on dog behaviour, and suggests that the predictive validity of behavioural assessments in dogs could be improved by systematically accounting for both inter- and intra-individual variation.

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

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

Authors' contributions. C.G. and R.C.N. conceptualized the study, revised the manuscript and wrote the final version. C.G. obtained the data, conducted the statistical analyses and drafted the initial manuscript.

Competing interests. We declare we have no competing interests.

Funding. C.G. and R.C.N. are employed by the Norwegian University of Life Sciences. No additional funding was required for this study.

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

References

- Bell AM, Hankison SJ, Laskowski KL. 2009 The repeatability of behaviour: a meta-analysis. *Anim. Behav.* **77**, 771–783. (doi:10.1016/j.anbehav.2008.12.022)
- Biro PA, Adriaenssens B, Cole AEBJ, Bronstein EJJ. 2013 Predictability as a personality trait: consistent differences in intraindividual behavioral variation. *Am. Nat.* **182**, 621–629. (doi:10.1086/673213)
- Dingemans NJ, Kazem AJN, Réale D, Wright J. 2010 Behavioural reaction norms: animal personality meets individual plasticity. *Trends Ecol. Evol.* **25**, 81–89. (doi:10.1016/j.tree.2009.07.013)
- Bridger D, Bonner SJ, Briffa M. 2015 Individual quality and personality: bolder males are less fecund in the hermit crab (*Pagurus bernhardus*). *Proc. R. Soc. B* **282**, 20142492. (doi:10.1098/rspb.2014.2492)
- Cleasby IR, Nakagawa S, Schielzeth H. 2015 Quantifying the predictability of behaviour: statistical approaches for the study of between-individual variation in the within-individual variance. *Methods Ecol. Evol.* **6**, 27–37. (doi:10.1111/2041-210X.12281)
- Mitchell DJ, Fanson BG, Beckmann C, Biro PA. 2016 Towards powerful experimental and statistical approaches to study intraindividual variability in labile traits. *R. Soc. open sci.* **3**, 160352. (doi:10.1098/rsos.160352)
- Stamps JA, Briffa M, Biro PA. 2012 Unpredictable animals: individual differences in intraindividual variability (IIV). *Anim. Behav.* **83**, 1325–1334. (doi:10.1016/j.anbehav.2012.02.017)
- Westneat DF, Wright J, Dingemans NJ. 2015 The biology hidden inside residual within-individual phenotypic variation. *Biol. Rev.* **90**, 729–743. (doi:10.1111/brv.12131)
- DeWitt TJ. 2016 Expanding the phenotypic plasticity paradigm to broader views of trait space and ecological function. *Curr. Zool.* **62**, 463–473. (doi:10.1093/cz/zow085)
- Scheiner SM. 2014 The genetics of phenotypic plasticity XIII. Interactions with developmental instability. *Ecol. Evol.* **4**, 1347–1360. (doi:10.1002/ece3.1039)
- Tonsor SJ, Elnacass TW, Scheiner SM. 2013 Developmental instability is genetically correlated with phenotypic plasticity, constraining heritability and fitness. *Evolution* **67**, 2923–2935. (doi:10.1111/evo.12175)
- Oates AC. 2011 What's all the noise about developmental stochasticity? *Development* **138**, 601–607. (doi:10.1242/dev.059923)
- Fratkin JL, Sinn DL, Patal EA, Gosling SD. 2013 Personality consistency in dogs: a meta-analysis. *PLoS ONE* **8**, e54907. (doi:10.1371/journal.pone.0054907)
- Wilsson E, Sundgren P-E. 1998 Behaviour test for eight-week old puppies—heritabilities of tested behaviour traits and its correspondence to later behaviour. *Appl. Anim. Behav. Sci.* **58**, 151–162. (doi:10.1016/S0168-1591(97)00093-2)
- Sinn DL, Gosling SD, Hilliard S. 2010 Personality and performance in military working dogs: reliability and predictive validity of behavioral tests. *Appl. Anim. Behav. Sci.* **127**, 51–65. (doi:10.1016/j.applanim.2010.08.007)
- Riemer S, Müller C, Virányi Z, Huber L, Range F. 2014 The predictive value of early behavioural assessments in pet dogs: a longitudinal study from neonates to adults. *PLoS ONE* **9**, e101237. (doi:10.1371/journal.pone.0101237)

17. Kis A, Klausz B, Persa E, Miklósi Á, Gácsi M. 2014 Timing and presence of an attachment person affect sensitivity of aggression tests in shelter dogs. *Vet. Rec.* **174**, 196. (doi:10.1136/vr.101955)
18. Serpell JA, Duffy DL. 2016 Aspects of juvenile and adolescent environment predict aggression and fear in 12-month-old guide dogs. *Front. Vet. Sci.* **3**, 49. (doi:10.3389/fvets.2016.00049)
19. Robinson LM, Thompson RS, Ha JC. 2016 Puppy temperament assessments predict breed and American Kennel Club group but not adult temperament. *J. Appl. Anim. Welf. Sci.* **19**, 101–114. (doi:10.1080/10888705.2015.1127665)
20. Marder AR, Shabelansky A, Patronek GJ, Dowling-Guyer S, D'Arpino SS. 2013 Food-related aggression in shelter dogs: a comparison of behavior identified by a behavior evaluation in the shelter and owner reports after adoption. *Appl. Anim. Behav. Sci.* **148**, 150–156. (doi:10.1016/j.applanim.2013.07.007)
21. Mohan-Gibbons H, Weiss E, Slater M. 2012 Preliminary investigation of food guarding behavior in shelter dogs in the United States. *Animals* **2**, 331–346. (doi:10.3390/ani2030331)
22. Mormenent KM, Coleman GJ, Toukhsati SR, Bennett PC. 2015 Evaluation of the predictive validity of the Behavioural Assessment for Re-homing K9's (B.A.R.K.) protocol and owner satisfaction with adopted dogs. *Appl. Anim. Behav. Sci.* **167**, 35–42. (doi:10.1016/j.applanim.2015.03.013)
23. Poulsen AH, Lisle AT, Phillips CJC. 2010 An evaluation of a behaviour assessment to determine the suitability of shelter dogs for rehoming. *Vet. Med. Int.* **2010**, e523781. (doi:10.4061/2010/523781)
24. Rayment DJ, Groef BD, Peters RA, Marston LC. 2015 Applied personality assessment in domestic dogs: limitations and caveats. *Appl. Anim. Behav. Sci.* **163**, 1–18. (doi:10.1016/j.applanim.2014.11.020)
25. Hennessy MB. 2013 Using hypothalamic-pituitary-adrenal measures for assessing and reducing the stress of dogs in shelters: a review. *Appl. Anim. Behav. Sci.* **149**, 1–12. (doi:10.1016/j.applanim.2013.09.004)
26. Protopopova A. 2016 Effects of sheltering on physiology, immune function, behavior, and the welfare of dogs. *Physiol. Behav.* **159**, 95–103. (doi:10.1016/j.physbeh.2016.03.020)
27. Stephen JM, Ledger RA. 2005 An audit of behavioral indicators of poor welfare in kennelled dogs in the United Kingdom. *J. Appl. Anim. Welf. Sci.* **8**, 79–95. (doi:10.1207/s15327604jaws0802_1)
28. Rooney NJ, Gaines SA, Bradshaw JWS. 2007 Behavioural and glucocorticoid responses of dogs (*Canis familiaris*) to kennelling: investigating mitigation of stress by prior habituation. *Physiol. Behav.* **92**, 847–854. (doi:10.1016/j.physbeh.2007.06.011)
29. Patronek GJ, Bradley J. 2016 No better than flipping a coin: reconsidering canine behavior evaluations in animal shelters. *J. Vet. Behav.* **15**, 66–77. (doi:10.1016/j.jveb.2016.08.001)
30. Protopopova A, Wynne CDL. 2014 Adopter-dog interactions at the shelter: behavioral and contextual predictors of adoption. *Appl. Anim. Behav. Sci.* **157**, 109–116. (doi:10.1016/j.applanim.2014.04.007)
31. Martin JGA, Pirottay E, Petellez MB, Blumstein DT. 2017 Genetic basis of between-individual and within-individual variance of docility. *J. Evol. Biol.* **30**, 796–805. (doi:10.1111/jeb.13048)
32. Kruschke J. 2014 *Doing Bayesian data analysis: a tutorial with R, JAGS, and Stan*. New York, NY: Academic Press.
33. McElreath R. 2015 *Statistical rethinking: a Bayesian course with examples in R and Stan*. Boca Raton, FL: CRC Press.
34. Voith VL *et al.* 2013 Comparison of visual and DNA breed identification of dogs and inter-observer reliability. *Am. J. Soc. Res.* **3**, 17–29. (doi:10.1080/10888700902956151)
35. Goold C, Newberry RC. 2017 Aggressiveness as a latent personality trait of domestic dogs: testing local independence and measurement invariance. *PLoS ONE* **12**, e0183595. (doi:10.1371/journal.pone.0183595)
36. Owczarczak-Garstecka SC, Burman OH. 2016 Can sleep and resting behaviours be used as indicators of welfare in shelter dogs (*Canis lupus familiaris*)? *PLoS ONE* **11**, e0163620. (doi:10.1371/journal.pone.0163620)
37. R Development Core Team. 2016 *R: a language and environment for statistical computing*. Vienna, Austria: R Development Core Team.
38. Tastle WJ, Wierman MJ. 2007 Consensus and dissent: a measure of ordinal dispersion. *Int. J. Approx. Reason* **45**, 531–545. (doi:10.1016/j.ijar.2006.06.024)
39. Ruedin D. 2016 agrmt: Calculate Agreement or Consensus in Ordered Rating Scales. R package version 1.40.4.
40. Liddell TM, Kruschke JK. 2015 Analyzing ordinal data: support for a Bayesian approach. *SSRN*. (doi:10.2139/ssrn.2692323)
41. Foulley J-L, Jaffrézic F. 2010 Modelling and estimating heterogeneous variances in threshold models for ordinal discrete data via Winbugs/Openbugs. *Comput. Methods Programs Biomed.* **97**, 19–27. (doi:10.1016/j.cmpb.2009.05.004)
42. Kizilkaya K, Tempelman RJ. 2005 A general approach to mixed effects modeling of residual variances in generalized linear mixed models. *Genet. Sel. Evol.* **37**, 31. (doi:10.1186/1297-9686-37-1-31)
43. Nakagawa S, Schielzeth H. 2010 Repeatability for Gaussian and non-Gaussian data: a practical guide for biologists. *Biol. Rev.* **85**, 935–956. (doi:10.1111/j.1469-185X.2010.00141.x)
44. Brommer JE. 2013 Variation in plasticity of personality traits implies that the ranking of personality measures changes between environmental contexts: calculating the cross-environmental correlation. *Behav. Ecol. Sociobiol.* **67**, 1709–1718. (doi:10.1007/s00265-013-1603-9)
45. Stan Development Team. 2016 Stan modeling language users guide and reference manual. Version 2.15.0.
46. Lewandowski D, Kurowiczka D, Joe H. 2009 Generating random correlation matrices based on vines and extended onion method. *J. Multivar. Anal.* **100**, 1989–2001. (doi:10.1016/j.jmva.2009.04.008)
47. Watanabe S. 2010 Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *J. Mach. Learn. Res.* **11**, 3571–3594.
48. Stan Development Team. 2016 Rstan: R Interface to Stan. R package version 2.14.1.
49. Hiby EF, Rooney NJ, Bradshaw JWS. 2006 Behavioural and physiological responses of dogs entering re-homing kennels. *Physiol. Behav.* **89**, 385–391. (doi:10.1016/j.physbeh.2006.07.012)
50. Svartberg K, Forkman B. 2002 Personality traits in the domestic dog (*Canis familiaris*). *Appl. Anim. Behav. Sci.* **79**, 133–155. (doi:10.1016/S0168-1591(02)00121-1)
51. Hsu Y, Serpell JA. 2003 Development and validation of a questionnaire for measuring behavior and temperament traits in pet dogs. *J. Am. Vet. Med. Assoc.* **223**, 1293–1300. (doi:10.2460/javma.2003.223.1293)
52. Jones AC, Gosling SD. 2005 Temperament and personality in dogs (*Canis familiaris*): a review and evaluation of past research. *Appl. Anim. Behav. Sci.* **95**, 1–53. (doi:10.1016/j.applanim.2005.04.008)
53. Riemer S, Müller C, Virányi Z, Huber L, Range F. 2016 Individual and group level trajectories of behavioural development in border collies. *Appl. Anim. Behav. Sci.* **180**, 78–86. (doi:10.1016/j.applanim.2016.04.021)
54. Cramer AOJ, Borkulo CD, Giltay EJ, Maas HLJ, Kendler KS, Scheffer M, Borsboom D. 2016 Major depression as a complex dynamic system. *PLoS ONE* **11**, e0167490. (doi:10.1371/journal.pone.0167490)
55. Wichers M *et al.* 2016 Critical slowing down as a personalized early warning signal for depression. *Psychother. Psychosom.* **85**, 114–116. (doi:10.1159/000441458)
56. David JT, Cervantes MC, Trosky KA, Salinas JA, Delville Y. 2004 A neural network underlying individual differences in emotion and aggression in male golden hamsters. *Neuroscience* **126**, 567–578. (doi:10.1016/j.neuroscience.2004.04.031)
57. Betini GS, Norris DR. 2012 The relationship between personality and plasticity in tree swallow aggression and the consequences for reproductive success. *Anim. Behav.* **83**, 137–143. (doi:10.1016/j.anbehav.2011.10.018)
58. Dewitt TJ, Sih A, Wilson DS. 1998 Costs and limits of phenotypic plasticity. *Trends Ecol. Evol.* **13**, 77–81. (doi:10.1016/S0169-5347(97)01274-3)
59. Highcock L, Carter AJ. 2014 Intraindividual variability of boldness is repeatable across contexts in a wild lizard. *PLoS ONE* **9**, e95179. (doi:10.1371/journal.pone.0095179)
60. Araya-Ajoy YG, Dingemans NJ. 2017 Repeatability, heritability, and age-dependence of seasonal plasticity in aggressiveness in a wild passerine bird. *J. Anim. Ecol.* **86**, 227–238. (doi:10.1111/1365-2656.12621)
61. Dykiert D, Der G, Starr JM, Deary IJ. 2012 Age differences in intra-individual variability in simple and choice reaction time: systematic review and meta-analysis. *PLoS ONE* **7**, e45759. (doi:10.1371/journal.pone.0045759)