

# Predicting individual shelter dog behaviour after adoption using longitudinal behavioural assessment: a hierarchical Bayesian approach

Conor Goold<sup>1,2</sup> and Ruth C. Newberry<sup>2</sup>

<sup>1</sup>School of Biology, Faculty of Biological Sciences, University of Leeds, UK, LS2 9JT

<sup>2</sup>Department of Animal and Aquacultural Sciences, Faculty of Biosciences, Norwegian University of Life Sciences, Ås, Norway

---

**Abstract.** Predicting the behaviour of shelter dogs after adoption is an important, but difficult, endeavour. A reliance on behavioural test batteries, considerable differences between shelter and post-adoption environments, between- and within-individual heterogeneity in behaviour, uncertainty in behavioural predictions, and measurement error all hinder the accurate assessment of future dog behaviour. This study integrated 1) a longitudinal behaviour assessment methodology with 2) a novel joint hierarchical Bayesian mixture model that accounted for individual variation, missing data and measurement error to predict the behaviour of dogs post-adoption. We analysed shelter behavioural observations (> 28,000 complete records in total) and post-adoption behavioural reports (from telephone surveys) across eight contexts from 241 dogs. Dog behaviour at the shelter correlated positively with behaviour post-adoption within contexts (0.38; 95% highest density interval: [0.20, 0.55]), and behaviour within contexts explained more of the behavioural variance than differences between dogs or between contexts alone. Although measurement error was higher post-adoption than at the shelter, we found few differences in individual-level, latent probabilities of different behaviours post-adoption versus at the shelter. This good predictive ability was aided particularly by accurate representation of uncertainty in individual-level behavioural predictions. We conclude that longitudinal assessments paired with a sufficient inferential framework to model latent behavioural profiles with uncertainty enable reasonably accurate assessments about post-adoption shelter dog behaviour.

---

# 1 Introduction

Prediction is a defining, yet challenging, goal of science [1], uniting a diverse range of disciplines including climate research, personality psychology, social policy, and complex systems science [?]. Likewise, accurately predicting the future behaviour of domestic dogs (*Canis familiaris*) is a challenge for professional organisations, dog breeders, and dog owners alike. In particular, animal shelters have the difficult task of assessing how dogs will behave in a variety of circumstances from limited behavioural information and making decisions about the suitability of those dogs for placement into new homes.

Animal shelters frequently employ standardised behaviour tests that evaluate specific behavioural or *personality* traits through reconstructions of situations relevant to life outside the shelter environment [?, ?, ?, ?, ?]. For example, to assess a dog’s level of aggressiveness around food (i.e. food guarding), shelter staff may present a dog with a bowl of food and record the dog’s response to a plastic hand approaching, touching or trying to remove the food bowl [?, ?, ?]. Such tests are usually conducted once, soon after arrival at the shelter, yielding an overall score used to help determine a dog’s suitability for adoption. However, criticisms have been raised about the ability for these protocols to provide accurate predictions of future dog behaviour. For example, the feasibility of carrying out a battery of such behaviour tests with high reliability and standardisation in the time-constrained shelter environment is limited given that test batteries often require at least an hour per dog [?, ?]; although see [?] for a shorter protocol). Moreover, a high number of false positives have been recorded for food guarding behaviour (i.e. displaying food aggression during shelter evaluations but not in the new home; see [?, ?] and has led to calls for abandoning standardised food tests given dog return rates to shelters do not appear to increase if the tests are not conducted [?]. Behaviour in other contexts might be similarly difficult to predict post-adoption [?, ?]; although see [?].

Ways to improve predictions about individual dog behaviour post-adoption are not entirely clear. While emphasis has been placed on improving the reliability and validity of shelter dog tests, progress is hampered by a lack of clarity of what those terms mean [?, ?], which is also a problem for the wider scientific community (e.g. see [?, ?, ?]. For instance, [?] demonstrate that even with high sensitivity and specificity, the probability of a dog showing aggression in the new home after a positive test in the shelter is unlikely to be higher than 50% and is probably closer to  $\sim 30\%$  (see also [?]). Tests reported to predict post-adoption behaviour successfully (e.g. [?, ?, ?] have been supported based on the statistical significance of linear associations (e.g. significant correlations or regression coefficients) between shelter and post-adoption behaviour (i.e. predictive ‘validity’; [?]. In contrast, arguments against the efficacy of behavioural tests highlight evidence of low predictive ‘ability’ for individual dogs (i.e. estimates of true/false positive and negative rates; [?].

As an alternative to standardised testing, several authors and organisations have empha-

sised the collection of daily behavioural observation records from dogs in shelters, in addition to other information (e.g. foster reports, pre-surrendering interviews). A holistic profile of each dog’s behaviour and welfare needs is then formulated and used to inform adoption decisions [?, ?, ?, ?, ?]. Such longitudinal approaches refrain from crudely labelling individual dogs as, for example, aggressive or non-aggressive based on a single behavioural testing outcome. However, research is still scarce on how best to implement this approach which requires summarising swathes of information generated on each dog in a manner to be of practical use for shelters.

We have previously reported on the behaviour of dogs while in shelters based on data collected using the longitudinal assessment methodology implemented by a large UK shelter organisation [?, ?]. The assessment relies on the spontaneous reporting of behavioural observations made by shelter staff during everyday contexts (e.g. walking past the dog’s kennel, putting a food bowl into the kennel, clipping on the lead or touching the collar). From these accumulated reports, [?] used the framework of behavioural reaction norms [?, ?] to partition individual variation in the behaviour of over 3,000 dogs when around unknown people (from almost 20,000 behavioural reports in total) into three components: personality (i.e. difference in average behaviour between individuals), plasticity (i.e. within-individual behavioural change over time) and predictability (i.e. within-individual residual variance). Accounting for all three components improved the estimated out-of-sample predictive accuracy of the statistical models.

While accounting for individual variation in behaviour across time at the shelter is important, making predictions about dog behaviour post-adoption based on longitudinal information about behaviour at the shelter brings its own challenges. Deciding on what statistical moments (e.g. location, scale) of the longitudinal shelter data should predict post-adoption behaviour is not obvious. Diagnostic statistics such as false positives or false negatives lose clarity because longitudinal assessments acknowledge that dog behaviour can change depending on time and context [?]. Inferential tools that can estimate the individual-level, latent probabilities of dogs showing different behaviours at different time points are required, rather than focusing on occurrences of single behaviours. The realities of collecting observational data in busy shelter environments also means that data may present significant patterns of non-ignorable missingness, substantially unbalanced repeated measures, and be subject to more measurement error than standardised assessments. Altogether, determining the predictive value of ‘on-going’ shelter assessments will require careful analysis of longitudinal data using statistical models that appropriately account for the data-generating processes.

The goal of the present study was to assess the correspondence between behaviour shown at the shelter, collected using the longitudinal assessment methodology, and the behaviour reported during telephone interviews post-adoption. Our secondary goal is to present a novel application of joint Bayesian hierarchical mixture modelling to compare predictions of dog behaviour recorded in shelters and post-adoption. Importantly, our approach accounts

Demographic variable	Mean (SD) or $n$ dogs
Number of observations per dog	140.7 (224.3)
Days at the shelter per dog	32.1 (41.4)
Estimated age at departure from shelter (years)	3.6 (2.8)
Weight (kg; last record at shelter)	16.5 (9.6)
Source (before surveys): gift/stray/return ( $n$ )	164/51/26
Sex: female/male ( $n$ )	102/139
Neutered: before arrival/at shelter/not/missing ( $n$ )	84/145/11/1
Rehoming centre: London/Old Windsor/Brands Hatch ( $n$ )	71/93/77

Table 1: Dog demographic variables collected at the shelter.

for i) individual variation in both average behaviour (i.e. personality) and behavioural change (i.e. plasticity; due to only two post-adoption interviews, predictability was not estimated) across shelter and post-adoption time periods, ii) missing data and their potential correlates, and iii) measurement error in the behavioural reports. We report on predictive validity by estimating correlations between personality and plasticity between shelter and post-adoption time periods. Predictive ability is examined by comparing model predictions for each dog’s average behaviour in the shelter and post-adoption to determine whether dogs exhibited credible changes in their behaviour. We present specific examples of individual-level behaviour in our results and discussion.

## 2 Materials & Methods

### 2.1 Subjects

Behavioural data were gathered retrospectively from new owners of 265 dogs adopted from Battersea Dogs and Cats Home, United Kingdom between May 2016 and May 2017. Details of the shelter environment can be found in [?] and [?]. Of the 265 dogs, only 241 dogs could be matched with behavioural and demographic data from the shelter’s database. Demographic information about these dogs is presented in Table 1. These dogs were generally of unknown heritage (i.e. not in a breed registry) and were derived from different kennels at three different shelter facilities. Demographic information of the new owners was collected but not analysed here due to inconsistencies in the responses. The original dog identification numbers assigned to dogs by the shelter were removed prior to analysis.

## 2.2 Data collection

### 2.2.1 Shelter behaviour

As previously described by [?] and [?], shelter employees observed dogs in a variety of naturally occurring contexts and recorded each dog’s behaviour using a context-specific list of mutually exclusive behavioural codes. The analysis presented here focussed on eight core ‘on-site’ contexts (Table 2). The *Interactions with dogs* context was a combination of the original *Interactions with female* and *male dogs* contexts, respectively, because the post-adoption data did not distinguish between interactions with male and female dogs (as adopters may not have been fully aware of the sex of other dogs that their dog met). Observation records were entered as often as possible given the frequency at which the context occurred and time constraints on shelter staff entering data. We considered any day that a dog did not receive an observation record in a context as a missing observation so that patterns of missingness could be modelled (see Table 2 for frequency of missing daily observation records). In total, there were 28,445 complete observation records and 61,952 missing and complete observation records for the 241 studied dogs in the shelter dataset.

There were between ten and sixteen possible behavioural codes depending on the context, arranged on a scale of perceived ease of adoption or desirability to adopters (see Supplementary Materials for the full list of behavioural codes). The longitudinal assessment further categorised the behavioural codes into green, amber and red categories: green codes indicated that the dog’s behaviour was suitable for immediate adoption and/or would be easy for all owners to manage once adopted (e.g. the dog was friendly or excited when meeting new people), amber codes indicated that some training or management might be needed to manage or improve behaviour (e.g. the dog was nervous when meeting new people and slow to meet with them), and red codes that the dog needed an individualised training program and improvement in behaviour to facilitate adoption or would be the most difficult for adopters to manage (e.g. the dog showed aggression when meeting unfamiliar people). A dog’s suitability for adoption was decided based on behaviour across all contexts and days after arrival.

The behavioural scores were analysed using the green, amber and red ordinal scale, rather than using the individual behavioural codes of each ethogram. This was chosen: i) to place the behaviour records across contexts on a comparable scale, allowing all the data to be analysed within the same statistical model; ii) to capture the main variation in the data given that some individual codes were seldom used; iii) because previous analysis indicated higher consensus among shelter staff when rating behaviour using the green, amber and red codes than the individual codes (see [?]; and iv) it was more practically relevant because the shelter was interested to find out if their assessments at the shelter differed greatly from a dog’s behaviour post-adoption (i.e. there was a change from green to red codes), but not necessarily if there was a change between codes within the same colour category. When more

Context	Description	Median (IQR)	% Missing (SD)
Handling (HND)	Informal handling by people (e.g. stroking non-sensitive areas, touching or holding the collar, fitting a harness or lead).	13 (14)	40 (20)
In kennel (KNL)	Inside the kennel	14 (14)	37 (21)
Outside the kennel (OKNL)	Outside the kennel	13 (14)	42 (19)
Interactions with familiar people (FPL)	Interacts with familiar people (interacted with at least once before) outside the kennel who approach, make eye contact, speak to or attempt to make physical contact with the dog.	13 (14)	44 (18)
Interactions with unfamiliar people (UFPL)	Interacts with familiar people (never interacted with before) outside the kennel who approach, make eye contact, speak to or attempt to make physical contact with the dog.	6 (5)	75 (15)
Eating food (FOOD)	Eating food from a container, toy or while being hand-fed	14 (13)	39 (20)
Interactions with toys (TOYS)	Interacting with dog toys	8 (13)	58 (25)
Interactions with dogs (DOGS)	Meeting dogs outside the kennel during structured interactions and/or spontaneous meetings	5 (4)	76 (12)

Table 2: *Observational contexts at the shelter with corresponding median and interquartile range (IQR; due to skewness) of the number of records per dog, and average percentage (standard deviation) of missing daily records per dog.*

than one record was made of the same dog in the same context on the same day, the most ‘severe’ code assigned was retained for analysis.

### 2.2.2 Post-adoption behaviour

On the point of adoption, adopters were asked to participate in a study evaluating the predictive accuracy of the shelter’s behavioural assessment and were given full details of the procedure (consent form provided in the Supplementary Materials). Consenting adopters each received two phone calls from a canine behaviourist at the shelter who surveyed them about their dog’s behaviour using a standardised set of questions (see Supplementary Materials). The majority of phone calls were made by one behaviourist although no records were kept on which behaviourist made the calls. The study plan was to record behaviour using telephone interviews at approximately 2-3 weeks and 5-6 weeks after adoption but the

actual days after adoption were more variable. The average number of days after adoption for the first phone call was 19.8 (sd = 4.9; range = 7 to 51), and 32.7 days (sd = 11.1; range = 14 to 72) for the second call, with an average of 20.7 (sd = 5.6; range = 4 to 42) days between first and second interviews for those dogs with two completed surveys. Only 150 dogs (62%) had two fully complete surveys. One dog was returned after two post-adoption surveys and was subsequently adopted and received two surveys from their second owner; we randomly chose the first set of surveys on this dog for inclusion in the analysis. A different dog had been returned after their first survey and was subsequently adopted and received two surveys from their second owner; we chose the second owners' surveys for this dog. Six of the dogs who only had a first survey completed had been returned to the shelter.

The first questions gathered information about the adopters' dog ownership experience and the post-adoption environment. The remaining questions (Table 3) enquired about how the dog reacted in situations comparable to the eight behavioural observation contexts in the shelter assessment (Table 2). The questions were used as a guide during the post-adoption interviews, with the behaviourists providing more detailed descriptions as needed. The *In kennel* and *Out of kennel* contexts at the shelter were transformed to *In house* and *Out of house* (where house referred to the residence in which the dog lived with the adopter). The questions were open-ended, with the behaviourist encouraging adopters to describe the dog's behaviour in each situation rather than label the behaviour. Subsequently, the behaviourist chose the behavioural code for each context that best described the behaviour. Adopters were allowed to respond 'No opinion', which was treated as a missing value. Phone calls were not recorded, precluding assessment of the reliability of behaviour coding during phone calls. Data were handled anonymously by the authors, who received only the dogs' identification numbers to enable matching of the post-adoption reports with the shelter records.

## 2.3 Statistical analysis

### 2.3.1 Theoretical approach

The green, amber and red codes from shelter and post-adoption time points were analysed using a joint or 'shared parameter' (e.g. [?, ?] Bayesian hierarchical model, which accounted for two data complexities. First, we treated the missing data present in the shelter records (Table 2) and post-adoption surveys (Table 3) as non-ignorable (i.e. not random) because they likely depended on other variables. For example, the number of people unfamiliar to a dog decreases with time spent at the shelter or time after adoption, leading to potentially more missing records in the *Interactions with unfamiliar people* context. A dog's overall behaviour could also impact missing values. For example, dogs with more green codes overall and who do not change their behaviour may not receive as many records due to a lack of priority relative to dogs who show more worrying or varying behaviour. The participation or attrition rate of adopters in the telephone surveys may also have depended on their dog's

Context (abbreviation)	Survey question	<i>n</i> (%) with 2 surveys
Handling (HND)	How has the dog behaved while being handled or restrained (informal, i.e. not in a veterinary context)?	149 (61.8)
In house (HOUSE)	How has the dog behaved in the house?	149 (61.8)
Outside of house (OTSD)	How has the dog behaved outside on walks?	147 (61.0)
Interactions with familiar people (FPL)	How has the dog behaved when meeting familiar people?	144 (59.8)
Interactions with unfamiliar people (UFPL)	How has the dog behaved when meeting unfamiliar people	149 (61.8)
Eating food (FOOD)	How has the dog behaved with their food?	150 (62.2)
Interactions with toys (TOYS)	How has the dog behaved with toys?	149 (61.8)
Interactions with dogs (DOGS)	How has the dog behaved when meeting another dog for the first time?	144 (59.8)

Table 3: Post-adoption survey behaviour questions and response rates (number and percentage of owners responding to the question over two surveys).

behaviour post-adoption.

A second complexity was the high proportion ( $\sim 90\%$ ) of green codes, which was found in previous analysis of the same shelter assessment [?]. Some proportion of the green codes may have been recorded due to other processes, such as staff members not observing the dog’s behaviour directly or forgetting how a dog behaved but inputting a green code anyway. Adopters may have also described their dog’s behaviour in terms consistent with green codes (i.e. not reporting any problems) during phone calls even if the dog’s behaviour might be more consistent with amber or red code behaviour. Thus, some green codes were potentially ‘inflated’, leading to a greater probability mass for green codes beyond that explained by the data generating processes of the behavioural scale alone. Similar assumptions are applied to count data with a high number of zero values (zero inflation; [?], and have also been used to understand a high probability mass for ‘I don’t know’ responses [?] and presumed face-saving ‘Neither agree nor disagree’ responses [?] on ordinal survey scales.

To account for the above complexities, we specified a custom mixture model for the probability of different codes ( $c = \text{missing, green, amber, red} = 0, 1, 2, 3$ ) for case  $i$  (either the days after arriving at the shelter or the days after adoption), dog  $j$  and context  $g$ :

$$p(y_{ijg}^c) = \begin{cases} \psi_{ijg} & \text{if } y_{ijg} = 0 \\ (1 - \psi_{ijg}[\kappa + (1 - \kappa)\pi_{ijg}^c]) & \text{if } y_{ijg} = 1 \\ (1 - \psi_{ijg}(1 - \kappa)\pi_{ijg}^c) & \text{if } 1 < y_{ijg} \leq 3 \end{cases} \quad (1)$$



The parameter  $\psi_{ijg}$  is a ‘hurdle’ probability of missing data for each dog and context (see [?] for a similar approach to handling missing data), which is modelled as a Bernoulli trial from the data:

$$y_{ijg}^c|_{c=0} \sim \text{Bernoulli}(\psi_{ijg}) \quad (2)$$

The complement,  $1 - \psi_{ijg}$ , is the probability of either a staff member choosing to input an observation in the shelter or an adopter participating in the telephone survey, respectively. The  $\kappa$  parameter is a mixing term that describes the probability of a green code being drawn from the inflation component. Consequently,  $1 - \kappa$  is the probability of a code being non-inflated. A non-missing and non-inflated code of category  $c$  (for a specific case, dog and context) occurs with probability  $(1 - \psi_{ijg})(1 - \kappa)\pi_{ijg}^c$ , where  $\pi_{ijg}^c$  is defined using a cumulative ordinal probit model:

$$y_{ijg}^c|_{c=\{1,2,3\}, 1-\kappa} \sim \text{Categorical}(\pi_{ijg}^c) \quad (3)$$

$$\pi_{ijg}^c = \Phi\left(\frac{\theta_c - \mu_{ijg}}{\sigma}\right) - \Phi\left(\frac{\theta_{c-1} - \mu_{ijg}}{\sigma}\right) \quad (4)$$

Equation 3 describes the non-missing and non-inflated code as categorically distributed with probability  $\pi_{ijg}^c$ , which is defined in equation 4 as the cumulative area under a latent standard normal distribution (where  $\Phi$  is the standard normal cumulative distribution function) between threshold parameters  $\theta_c$  and  $\theta_{c-1}$  ( $\boldsymbol{\theta} = \theta_0, \theta_1, \theta_2, \theta_3$ ). The probability of the first and last threshold parameters (i.e.  $\Phi(\theta_0)$ ,  $\Phi(\theta_3)$ ) were set to 0 and 1, respectively. To estimate the mean ( $\mu_{ijk}$ ) and standard deviation ( $\sigma$ ) on the scale of the ordinal data, we fixed  $\theta_1 = 1.5$  and  $\theta_2 = 2.5$  (see [?]).

### 2.3.2 Joint modelling

The data from the shelter and post-adoption time periods were each analysed by the model described above, which we describe as two distinct ‘sub-models’, both predicting the probability of missing data ( $\psi_{ijg}$  in equation 2 using a logit link function) and the predicted mean of the latent behavioural scale scores ( $\mu_{ijg}$  in equation 4 using the identity link function). Each sub-model included the dog-, context-, and dog  $\times$  context-varying intercept and slope parameters (for the number of days after arrival at the shelter or days after adoption, respectively). The joint structure of our model derives from correlations between these parameters across the shelter and post-adoption sub-models. For the ordinal probit regressions (equation 4, we included random intercepts and slope terms for dogs, contexts and the interaction between dogs and contexts (1,928 unique dog  $\times$  contexts combinations). Dog- and context-varying effects capture additive variation due to dogs and contexts, respectively, while the dog  $\times$  context interaction describes the non-additive, unique variation attributed to specific

dogs in specific contexts. For the missing data regressions, we included only random inter-  
 cepts. For each random effect term (dogs, contexts, and dogs  $\times$  contexts), this led to a 6 x 6  
 covariance matrix capturing the relationship between random intercepts and slopes for the  
 shelter and post-adoption behaviours.

### 2.3.3 Predictor variables

For the shelter sub-model, we included days after arrival at the shelter as an observation-level  
 predictor and the variables from Table 1 as dog-level predictors, excluding the total number  
 of observations (which was nearly perfectly collinear with length of stay) and rehoming  
 centre (which was not of interest *a priori*). The post-adoption sub-model included days after  
 adoption as an observation-level predictor and the same shelter demographic variables as in  
 the shelter sub-model as retrospective dog-level predictors, as well as the number of post-  
 adoption surveys about each dog. Sum-to-zero coding was used for categorical predictors.  
 Dog-level metric predictors were mean-centred and standardised by 1 standard deviation.  
 The number of days after arrival at the shelter was centered around the average length of  
 stay at the shelter (32.1 days; Table 1 and scaled by 1 standard deviation (73 days). The  
 number of days after adoption was centered around the average of each dog’s mean number  
 of days after adoption (26.3 days) and scaled by its standard deviation ( $\sim 12$  days). One  
 dog had missing neuter status data, two dogs were missing weight data and two dogs were  
 missing age at departure data. Due to the small amount of missing data, we imputed the  
 missing neuter status data point with the most frequent category (neutered on site; Table  
 1, and mean imputed the missing weight and age data points.

### 2.3.4 Estimation & inference

All data cleaning and post-processing was conducted in R version 3.6.1 [?]. We fit our model  
 using Bayesian estimation in the Stan programming language [?] using the terminal interface  
 CmdStan version 2.23.0 [?]. Stan employs Hamiltonian Monte Carlo, a type of Markov chain  
 Monte Carlo (MCMC) sampling algorithm, to sample from the posterior distribution. After  
 ensuring adequate mixing of multiple chains from the model ran on smaller sub-sets of  
 the data, we took 10,000 iterations from the posterior distribution using 2 MCMC chains  
 (1000 warmup and 5000 sampling). We summarise parameters by their means and 95%  
 highest density intervals (HDIs). For the ordinal data, we make predictions using both the  
 latent metric scale (parameters denoted by  $\beta$ ) and the corresponding posterior probabilities  
 of each ordinal category (probabilities denoted by  $\pi$ ). When evaluating population-level  
 parameters, we highlight discrepancies, where appropriate, between the estimates for average  
 dogs and contexts, and estimates that incorporate uncertainty due to variation across dogs  
 and contexts. The former represent estimates made at the mean of the random effects, while  
 the latter are made by marginalising over the random-effect distributions and are particularly

important for making accurate future predictions in hierarchical models (e.g. [?, ?]).

## 2.4 Ethical statement

Approval for the processing of personal data was obtained from the Norwegian Social Science Data Services (approval number 47080). Names and addresses of the participating new owners and their dogs were anonymised before data were passed to the authors.

## 2.5 Data & code accessibility

We provide complete mathematical details of the above model, supplementary files, the data, the Stan code, and the R code to reproduce the results reported here at <https://github.com/cmgoold/goold-newberry-lba>. Due to the non-standard statistical model, we also provide R and Stan code to simulate data and fit the model to recover parameter values for model validation at the same repository.

# 3 Results

## 3.1 Average behaviour

Green codes accounted for more than 95% of the observations in the raw shelter and post-adoption datasets, and for an average dog in an average context (i.e. at the mean of the random effects, and where all dog-level predictors equalled 0), there were no credible differences between the probability of green, amber and red codes at the shelter (at approximately 30 days after arrival) and post-adoption (at approximately 30 days post-adoption; probability differences:  $\pi_{\text{green}} = 0.00$ , 95% HDI: [-0.05, 0.06];  $\pi_{\text{amber}} = 0.00$ , 95% HDI: [-0.06, 0.05];  $\pi_{\text{red}} = 0.00$ , 95% HDI: [0.00, 0.00]). Similarly, for an average dog in an average context, behaviour tended to improve with every additional week at the shelter ( $\beta = -0.08$ ; 95% HDI: [-0.14, -0.02]), although due to ceiling effects the practical increase in green codes was minimal (from around 98% on arrival day to 99% on day 30). There was no clear relationship of days after adoption on post-adoption behaviour ( $\beta = 0.01$ ; 95% HDI: [-0.21, 0.24]). A dog's average probability of missing data at the shelter was  $\pi_{\text{missing}} = 0.54$  (95% HDI: [0.37, 0.73]), and the odds of missing codes increased with every additional week after arrival (odds: 1.04; 95% HDI: [1.03, 1.05]). The odds of missing data post-adoption increased sharply with every additional week after adoption (odds: 11.62; 95% HDI: [8.90, 14.00]), although the baseline probabilities of missing data at 30 days post-adoption between dogs with one survey were  $\pi_{\text{missing}} = 0.45$  (95% HDI: [0.26, 0.65]) compared to essentially 0 probability of missing data for dogs with two surveys, highlighting that within-survey missing data was highly infrequent.

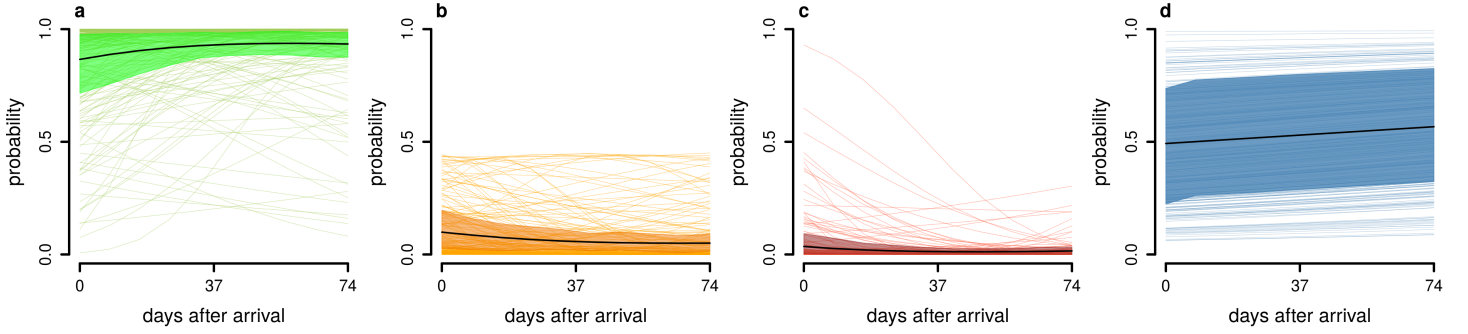


Figure 1: Probability of green (a), amber (b), red (c) and missing (d) codes across days after arrival at the shelter (average length of stay + 1 standard deviation) accounting for variation in behaviour across dogs, contexts and dog  $\times$  context. Thick black lines show posterior mean estimates, coloured bands show the 95% highest density interval of the mean, and thin coloured lines ( $n = 241$ ) show one random sample from one randomly-chosen context from each dog’s posterior distribution.

The probability of green-code inflation at the shelter was  $\kappa = 0.04$  (95% HDI: [0.00, 0.09]) at the shelter, and  $\kappa = 0.20$  (95% HDI: [0.03, 0.36]) post-adoption.

### 3.2 Variation across contexts

330 Incorporating variation across dogs and contexts demonstrated substantially more uncertainty in the model estimates both at the shelter (Figure 1 and post-adoption (Figure 2. At the shelter, the *Interactions with dogs* context had the highest probability of red codes ( $\pi_{red} = 0.04$ ; 95% HDI: [0.02, 0.06]; Figure S1a), and dogs changed their behaviour the most in the *Handling* context ( $\beta = -0.20$  95% HDI: [-0.28, -0.14]; Figure S1b). Post-adoption, 335 the *Eating food* context had the highest probability of red codes ( $\pi_{red} = 0.01$ ; 95% HDI: [0.00, 0.05]; Figure S2a), although there were no strong differences between contexts in the average amount of behavioural change across days after adoption (Figure S2b). Missing codes were particularly prominent at the shelter for the *Interactions with unfamiliar people* context ( $\pi_{missing} = 0.81$ ; 95% HDI: [0.74, 0.89]).

### 340 3.3 Interaction between dogs and contexts

Both at the shelter and post-adoption, behaviour of dogs within contexts (the dog  $\times$  context-varying intercepts) varied more than between dogs across all contexts or between contexts across all dogs. This is illustrated by higher repeatability (i.e. proportion of total variance explained; see Supplementary Materials for calculation) for the dog  $\times$  context-varying intercept 345 parameters compared to the additive effects of dogs or contexts alone (Figure 3a). Repeatability post-adoption was approximately 20% higher (but more uncertain) post-adoption for dogs within contexts than at the shelter ( $R_{diff} = -0.23$ ; 95% HDI: [-0.41, -0.05]) because

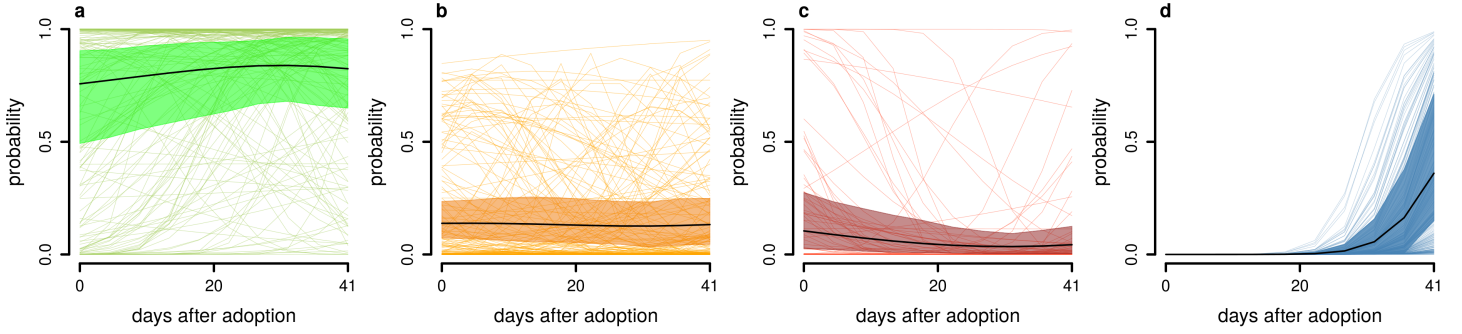


Figure 2: Probability of green (a), amber (b), red (c) and missing (d) codes across days after adoption (average days after adoption interview + 1 standard deviation) accounting for variation across dogs, contexts and dog  $\times$  context behaviour. Thick black lines show posterior mean estimates, coloured bands show the 95% highest density interval of the mean, and thin coloured lines ( $n = 241$ ) show one random sample from one randomly-chosen context from each dog’s posterior distribution.

the standard deviation of the dog  $\times$  context-varying intercepts was 1.40 (95% HDI: [1.10, 1.72]) times higher post-adoption than at the shelter. The standard deviation of the dog  $\times$  context slope parameters was considerably larger at the shelter (describing behaviour over days after arrival) than post-adoption (describing behaviour over days after adoption), but had substantial uncertainty (ratio of standard deviations: 45.54; 95% HDI: [1.61, 96.43]).

### 3.4 Random effect correlations

The dog  $\times$  context-varying intercept parameters at the shelter and post-adoption had a moderate positive correlation ( $\rho = 0.38$ ; 95% HDI: [0.20, 0.55]; Figure 3b), whereas the correlations between dog-varying and context-varying intercepts were largely positive but their 95% HDIs spanned zero (Table S1). Correlations between the random-slope parameters at the shelter and post-adoption included zero. At the shelter, the dog  $\times$  context intercept and slope parameters had a positive correlation ( $\rho = 0.26$ ; 95% HDI: [0.02, 0.47]), meaning that dogs with higher values on the latent behavioural scale (higher changes of red codes) changed their behaviour less over time or, in some cases, were more likely to exhibit amber and red codes over time. Dogs with more positive slope parameters across contexts at the shelter had slightly higher chances of missing data at the shelter ( $\rho = 0.28$ ; 95% HDI: [0.01, 0.54]). Dogs with fewer missing data on average post-adoption also had larger intercept parameters ( $\rho = -0.38$ ; 95% HDI: [-0.71, -0.04]) and more positive slopes ( $\rho = -0.36$ ; 95% HDI: [-0.69, -0.02]) at the shelter, suggesting that adopters of dogs with overall higher chances of amber and red codes at the shelter had better response rates in the post-adoption phone calls. We note that average post-adoption behaviour had a mostly positive correlation with the amount of missing data post-adoption, although the 95% HDI was wide and just included zero ( $\rho = 0.25$ ; 95% HDI: [-0.07, 0.53]).

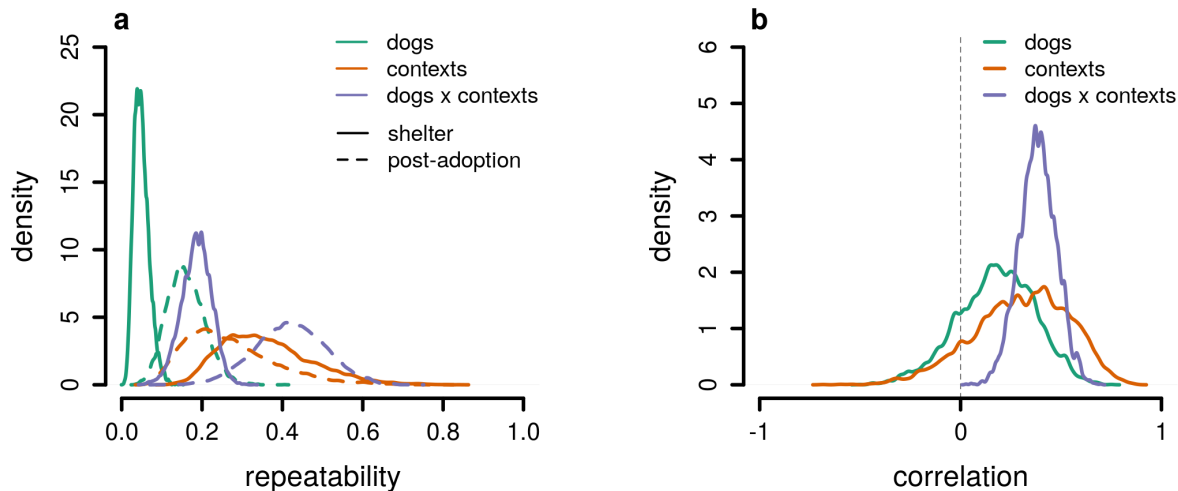


Figure 3: Proportions of total behavioural variance explained (repeatability) by dog-, context- and dog  $\times$  context-varying intercepts for shelter and post-adoption time periods (a), and correlations for dog-, context- and dog  $\times$  context-varying intercepts between shelter and post-adoption time periods (b; Table S1). Each parameter is evaluated at approximately 30 days after arrival at the shelter and 30 days after adoption.

### 3.5 Individual-level prediction of post-adoption behaviour

Behaviour was not static across shelter and post-adoption time points. A median of 48.50 (inter-quartile range: 32.75) and 39 (inter-quartile range: 66.75) dogs who received at least one amber and red code, respectively, within contexts at the shelter did not show amber or red codes post-adoption for the same contexts. Similarly, a median of 29.50 (inter-quartile range: 56.75) and 13 (inter-quartile range: 65.25) dogs who only showed green or amber codes at the shelter within contexts, received red codes in the post-adoption reports. As mentioned earlier, these diagnostic values only provide a crude estimate of behavioural differences between the shelter and post-adoption time periods, ignoring unbalanced repeated measures, measurement error and missing data. They further ignore each dog's latent probabilities of showing green, amber and red behaviours.

Therefore, we compared the posterior probabilities of green, amber and red codes on each dog's mean day after arrival at the shelter with the same codes' predicted probabilities on the day of the first and second telephone survey (where dogs with only the first survey had their second survey predictions estimated by the model). We summarised the resulting differences ( $\delta$ ) in code probabilities by counting the number and proportion of dogs whose 95% HDIs of the differences did not contain zero (Figure 4). In general, the number/proportion of dogs who had significantly different probabilities of green, amber or red codes post-adoption compared to the shelter was low (mean: 2.83; median: 0; sd: 7.51; range: 0 to 53). There were fewer overall differences between behaviour reported at the dogs' second surveys and the shelter

records than between the first survey and shelter records, likely due to fewer second surveys available resulting in increased uncertainty in the predictions. Most notably, approximately 10% of dogs ( $n = 24$ ) had lower probabilities of green codes, and 10% ( $n = 24$ ) and 13% ( $n = 31$ ) of dogs had higher chances of amber and red codes, respectively, in the *In kennel* context at the shelter than the corresponding *In house* context at the first survey, although this discrepancy had largely disappeared by the second survey phone call. Around 20% of dogs ( $n = 53$ ) had higher probabilities of green codes in the *Eating food* context at the shelter compared to the first survey, dropping to approximately 11% ( $n = 26$ ) at the second survey. A smaller number of dogs (9 and 5 at first and second surveys, respectively) had increased chances of red codes post-adoption in the *Eating food* context.

### 3.6 Dog-level predictors

Length of stay at the shelter had a positive (higher latent scale scores) but weak relationship with dogs' average behaviour at the shelter ( $\beta = 0.02$ ; 95% HDI: [0.01, 0.03]) and post-adoption ( $\beta = 0.02$ ; 95% HDI: [0.00, 0.03]). Dogs neutered before arrival had higher latent scale scores on average than intact dogs ( $\beta_{\text{no--yes}} = -0.38$ ; 95% HDI: [-0.73, -0.07]). Weight had a negative relationship (green codes increased with weight) with average behaviour post-adoption ( $\beta = -0.15$ , 95% HDI: [-0.26, -0.05]).

## 4 Discussion

This study has provided the first comprehensive analysis for predicting shelter dog behaviour post-adoption from a longitudinal behavioural assessment. We applied a novel Bayesian hierarchical model to account for both shelter and post-adoption behaviour, missing data, and measurement error within a single statistical framework, which allowed making individual-level inferences of how dogs behaved in specific contexts at the shelter and in their new homes. We found few differences between how likely dogs were to show green, amber and red-coded behaviour on average at the shelter, and their corresponding chances of showing the same behavioural codes in the post-adoption surveys. Moreover, in contrast to some previous studies reporting high numbers of behavioural problems in shelter dogs (e.g. [?], the behavioural reports in this study were largely dominated by green codes (i.e. the most favourable behaviours). This does not imply that dogs behaved exactly the same post-adoption as they did in the shelter, however. Due to the varying amounts of data available per dog per context, there was considerable uncertainty in individual-level, latent probabilities of green, amber and red codes, particularly in the post-adoption reports where dogs had a maximum of two records per context. This makes it plausible, for example, that a dog who only displayed green or amber codes at the shelter but showed red code behaviour post-adoption to have statistically similar latent probabilities of different behaviours at the shelter

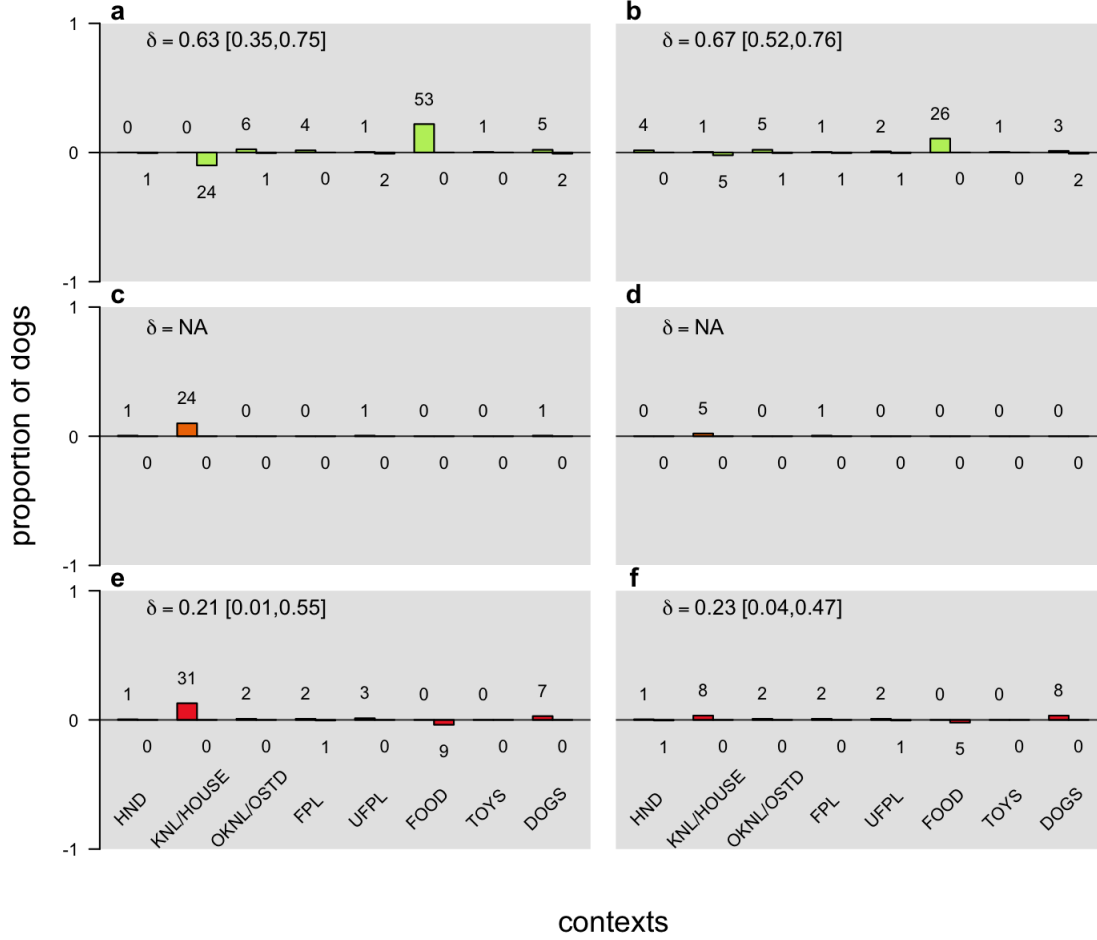


Figure 4: The proportion (bars) and number of dogs (text above each bar) who had significantly different probabilities of green (a-b), amber (c-d), and red (e-f) codes within contexts at the shelter compared to behaviour reported at the first (a, c, e) and second (b, d, f) post-adoption phone calls. Positive bars illustrate the proportion of dogs with higher code probabilities at the shelter compared to post-adoption; negative bars illustrate the proportion of dogs with lower probabilities at the shelter compared to post-adoption. A dog was included in the count if the 95% highest density interval of the differences in code probabilities for a dog did not include zero. The mean ( $\delta = [0, 1]$ ) and 95% highest density interval absolute probability difference is shown above each plot, except for panels c and d where there were too few dogs for its calculation.

and post-adoption. Far from being categorised as a ‘false negative’, longitudinal assessment and our inferential framework provides a richer, probabilistic assessment of individual dog behaviour on which to make more holistic assessments.

The hierarchical structure of our data allowed comparing variation in the behavioural records across dogs, across contexts, and across specific dogs  $\times$  contexts combinations. Both at the shelter and post-adoption, the dog  $\times$  context interaction explained the largest portion of behavioural variance (Figure 3a), and the correlations between shelter and post-adoption behaviour were only credibly larger than zero for the dog  $\times$  context-varying intercept parameters (Figure 3b). Neither how dogs behaved on average across contexts, nor behaviour



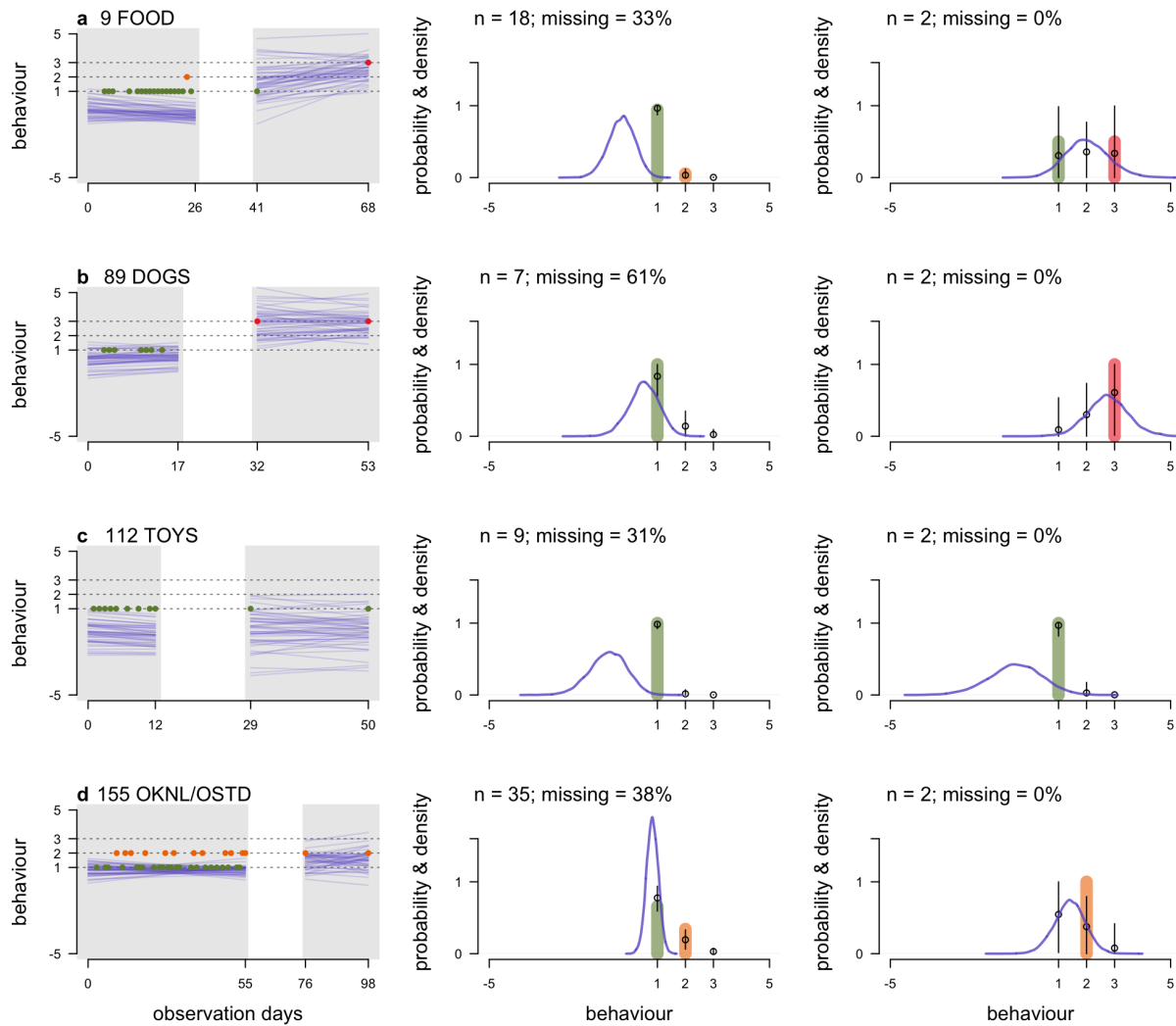


Figure 5: Posterior predictions for 4 specific dogs and contexts (a-d; anonymous dog ID number and context abbreviation shown above first panel of each plot; see Supplementary Materials for more examples). Left panels show behaviour over days after arrival at the shelter (first gray area) and behaviour over days after adoption (second gray area), where the coloured dots denote the raw ordinal data and blue lines 50 samples from the posterior distribution on the latent behavioural scale. Middle and right panels show the raw probability of green, amber and red codes (coloured bars) for the shelter (middle) and post-adoption (right), with the number of observations and percent missing data displayed above. Black dots and vertical line segments denote the mean and 95% highest density interval of predicted probabilities. Blue density curves show the underlying latent scale generating the ordinal codes for each dog.

within contexts irrespective of particular dogs, were particularly illuminating for predicting behaviour post-adoption. Repeatability was also around 20% higher post-adoption than in the shelter for the dog  $\times$  context behaviour, driven by greater heterogeneity in the post-adoption behaviour estimates than the shelter reports.

The positive correlation for dog behaviour within contexts at the shelter and post-

adoption informs us about their linear association, but the practical implications of this positive correlation are not easy to discern. One option is to compute rates of false positives and negatives (e.g. [?]) or positive and negative predictive values (e.g. [?]) to understand an assesment’s predictive ‘ability’ [?]. This medical diagnostic approach lends itself easily to scenarios whereby behavioural testing is conducted once at the shelter and once post-adoption, and for test protocols where the outcome measure is a binary pass or fail result. However, moving away from battery-style behavioural assessments in shelters necessitates moving away from labelling dogs as either problematic or not problematic, aggressive or non-aggressive, from limited behavioural information. Such sharp behavioural categorisations are at odds with the rationale for longitudinal assessments (e.g. [?, ?]) that encourages forming a holistic profile of a dog’s behaviour over different contexts, rather than focusing on any one single behaviour. Instead, we leveraged the posterior predictions of our model to determine how similar dog behaviour was across the post-adoption behaviour reports to each dog’s behaviour on average at the shelter. Hierarchical models are excellently suited to this task because they balance the information from each dog with what is typical across the whole sample. Thus, any one dog’s individual-level predictions are influenced by how many times it performed certain behaviours, the number of missing records for that dog, and the consistency of that dog’s behaviour in a specific context. This ‘partial pooling’ of information across dogs and contexts makes better predictions than non-hierarchical models or raw data because it balances over- or under-fitting to the data (e.g. [?, ?]).

To pick a few representative examples (see Supplementary Materials for more examples), Figure 5 illustrates the posterior predictions for four dogs, each in a different context, with the raw data for comparison. On the latent scale generating the ordinal data, dogs’ predictions for post-adoption behaviour (indicated by the sample of regression lines in the left-hand plots) tend to be more uncertain than their in-shelter predictions, where more data is available overall (e.g. compare shelter and post-adoption results in Figure 5a). The impact of the amount of data available can also be seen by inspecting the overall probabilities of different code colours in the middle and right panels of each plot, where dogs with fewer data points have more uncertainty (e.g. the 95% HDIs in the middle panel of 5b). What occurs on average in certain contexts also affects the individual-level predictions. For instance, the dog in Figure 5b received no red codes at the shelter in the *Interactions with dogs* context, but received red codes post-adoption. The overall increase in this dog’s probability of red codes in the latter context compared to the shelter is 56%, but with a 95% highest density interval that spans from 100% to -1%, meaning there is a small chance that red codes were just as, or more, likely at the shelter than post-adoption. While this might seem surprising, it occurs because red codes were most probable at the shelter in the *Interactions with dogs* context overall, so even dogs that show no red codes in this context have an increased chance of red codes compared to other contexts. The *Interactions with dogs* context also had one of the highest amounts of missing data (dog 89 in Figure 5b has only 7 observations and 61% of

missing data), and thus the predictions are more uncertain overall. By contrast, green codes  
480 were particularly dominant in the *Interactions with toys* context, and the amount of missing  
data was lower than average (Figure S1b), so dogs showing green codes in this context at  
the shelter and post-adoption have much tighter posterior estimates (e.g. see the predicted  
probabilities for the dog in Figure 5c).

For most contexts, the number of dogs whose 95% most likely values of the differences  
485 in probabilities of green, amber and red codes at the shelter compared to post-adoption did  
not include zero was low (Figure 4. This supports the ability for the longitudinal assessment  
to predict post-adoption behaviour, although we emphasise that this good predictive ability  
is in part due to the appropriate representation of uncertainty in the estimates, and the  
choice of 95% highest density intervals to summarise the posterior distribution. The greatest  
490 discrepancies between shelter and post-adoption behaviour appeared in the *In kennel/In  
house* and *Eating food* contexts. Between 10 and 15% of dogs in the former context were  
predicted to show more green codes, and fewer amber and red codes, at the first post-adoption  
phone call than at the shelter (Figure 4. This is perhaps unsurprising as the *In kennel* and *In  
house* contexts were the least comparable contexts across shelter and post-adoption periods,  
495 and a kennel would be expected to be a less stressful environment than the adoptive home.  
In contrast, around 20% of dogs at the first post-adoption survey, and around 11% at the  
second survey, showed lower chances of green codes than at the shelter in how they behaved  
around their food. The latter result is partly consistent with previous research questioning  
whether reasonable predictions of how dogs behave around food in the adopted home can  
500 be made from their behaviour at the shelter (e.g. [?, ?]. We also found that, post-adoption,  
the probability of red codes was highest in the *Eating food* context. However, the latter  
probability was still very low ( $\sim 1\%$ ), and apart from nine dogs at the first survey and  
five dogs at the second, we found no evidence for increased probabilities of red codes post-  
adoption for behaviour around the food bowl. Thus, observational data collected on how  
505 dogs behave around their food bowl may provide a more valid (in the predictive ability sense)  
indicator of post-adoption behaviour, even in this historically difficult-to-infer context.

Longitudinal behavioural assessments and adopter-reported behaviour are at risk of both  
missing data and measurement error, which are unlikely to be avoided entirely. Instead,  
understanding when and how they emerge will be key to accounting for any deleterious  
510 effects they have on behavioural data. We found that dogs' odds of missing data at the  
shelter were positively correlated with the amount of behavioural change across days after  
arrival. Given that the average coefficient describing behavioural change at the shelter was  
negative (indicating an improvement in behaviour through time), and that most dogs had  
largely negative individual slope estimates (i.e. deterioration in behaviour was rare), the  
515 latter correlation means that dogs who were less plastic were more likely to receive missing  
data on average at the shelter. Importantly, there was no evidence to suggest that dogs  
with less favourable behaviour at the shelter on average were more likely to receive missing

data (Table S1). We further found that owner response rates in the post-adoption surveys correlated more clearly with dog behaviour reported in the shelter than behaviour reported post-adoption: dogs with higher latent scale averages at the shelter, as well as dogs who changed less over time at the shelter, had lower chances of missing data post-adoption. We refrain from hypothesising-after-the-results-are-known (i.e. HARKing; [?]) on the causal mechanisms driving these correlations, but do emphasise the importance of understanding how response patterns (from shelter staff and new owners) are impacted by dog behaviour for the efficacy of longitudinal behaviour assessments.

The post-adoption behaviour reports had larger measurement error estimates than the shelter reports. Like estimates of zero-inflation in count data [?], we estimated a ‘green code inflation’, that is, the proportion of green codes reported due to processes different to the main process generating the behavioural data (i.e. the dog’s behaviour). In reality, the inflation component is a place-holder for more complex mechanisms, which could include shelter staff not observing a dog’s behaviour but inputting a green code anyway, or adopters reporting behaviours consistent with green codes, even when the real behaviour was less favourable. Green codes had approximately a 20% chance of being drawn from the inflation component of the model post-adoption, although there was large uncertainty (95% HDI: [3, 36]). By comparison, green code inflation was only 4% at the shelter. While we cannot disentangle the adopters’ own behavioural reports from how the shelter behaviourists interpreted the behavioural descriptions from the adopters, this implies that systematic error in the adopter-reported behavioural reports is important to quantify further. Unfortunately, no previous studies have assessed patterns of measurement error in shelter or post-adoption behaviour.

The modelling approach in this study is provided completely open source, and we encourage future studies on predicting shelter dog behaviour to consider principled statistical frameworks such as hierarchical Bayesian methods. For example, if more data was available on dogs post-adoption, this study could have made predictions of a sub-sample of shelter dogs’ post adoption behaviour reports left out of joint shelter–post adoption analysis. Moreover, due to the cumulative nature of Bayesian inference [?], the same modelling approach could be applied in real-time at shelters to provide predictions of shelter dog behaviour. The relative probabilities of dogs showing different types of behaviours could be integrated with qualitative information on each individual dog (e.g. pre-surrendering surveys; [?]) and expert knowledge to form a informed adoption decision and find a suitable new home.

## 5 Conclusion

Predicting the future behaviour of shelter dogs is a difficult task. Longitudinal assessments have been proposed as an alternative to battery-style, standardised testing, but research is still scarce on their implementation and efficacy. We have provided the first comprehensive

555 analysis of a longitudinal assessment to predict the behaviour of shelter dogs post-adoption  
using a novel application of Bayesian hierarchical modelling. While dog behaviour was not  
static, few dogs showed marked changes in their latent behavioural profiles at the shelter and  
post-adoption, supporting the use of longitudinal assessments and our inferential framework  
to inform adoption decisions alongside additional information and expert judgement. We  
560 have, further, demonstrated links between dog behaviour, missing data and measurement  
error, which should be accounted for in future analyses of longitudinal assessment data.

## References

- [1] Hofstadter, A. Explanation and necessity. *Philosophy and Phenomenological Research* **11**, 339–347 (1951).
- 565 [2] Sarewitz, D. & Pielke, R. Prediction in science and policy. *Technology in Society* **21**, 121–133 (1999).
- [3] van der Borg, J. A., Netto, W. J. & Planta, D. J. Behavioural testing of dogs in animal shelters to predict problem behaviour. *Applied Animal Behaviour Science* **32**, 237–251 (1991).
- 570 [4] Marston, L. C. & Bennett, P. C. Reforging the bond—towards successful canine adoption. *Applied Animal Behaviour Science* **83**, 227–245 (2003).
- [5] Mornement, K. M., Coleman, G. J., Toukhsati, S. & Bennett, P. C. A review of behavioral assessment protocols used by australian animal shelters to determine the adoption suitability of dogs. *Journal of Applied Animal Welfare Science* **13**, 314–329  
575 (2010).
- [6] Taylor, K. D. & Mills, D. S. The development and assessment of temperament tests for adult companion dogs. *Journal of Veterinary Behavior* **1**, 94–108 (2006).
- [7] Rayment, D. J., De Groef, B., Peters, R. A. & Marston, L. C. Applied personality assessment in domestic dogs: Limitations and caveats. *Applied Animal Behaviour Science*  
580 **163**, 1–18 (2015).
- [8] Mohan-Gibbons, H., Weiss, E. & Slater, M. Preliminary investigation of food guarding behavior in shelter dogs in the united states. *Animals* **2**, 331–346 (2012).
- [9] Mohan-Gibbons, H. *et al.* The impact of excluding food guarding from a standardized behavioral canine assessment in animal shelters. *Animals* **8**, 27 (2018).
- 585 [10] Marder, A. R., Shabelansky, A., Patronek, G. J., Dowling-Guyer, S. & D’Arpino, S. S. Food-related aggression in shelter dogs: A comparison of behavior identified by a behavior evaluation in the shelter and owner reports after adoption. *Applied animal behaviour science* **148**, 150–156 (2013).
- [11] Mornement, K., Toukhsati, S., Coleman, G. & Bennett, P. Reliability, validity and  
590 feasibility of existing tests of canine behaviour. In *AIAM Annual Conference on Urban Animal Management, Proceedings*, 11–18 (2009).
- [12] Poulsen, A., Lisle, A. & Phillips, C. An evaluation of a behaviour assessment to determine the suitability of shelter dogs for rehoming. *Veterinary medicine international* **2010** (2010).

- 595 [13] Christensen, E., Scarlett, J., Campagna, M., Houpt, K. A. *et al.* Aggressive behavior in adopted dogs that passed a temperament test. *Applied Animal Behaviour Science* **106**, 85–95 (2007).
- [14] Mornement, K. M., Coleman, G. J., Toukhsati, S. R. & Bennett, P. C. Evaluation of the predictive validity of the behavioural assessment for re-homing k9’s (bark) protocol and  
600 owner satisfaction with adopted dogs. *Applied Animal Behaviour Science* **167**, 35–42 (2015).
- [15] Bollen, K. S. & Horowitz, J. Behavioral evaluation and demographic information in the assessment of aggressiveness in shelter dogs. *Applied Animal Behaviour Science* **112**, 120–135 (2008).
- 605 [16] Patronek, G. J., Bradley, J. & Arps, E. What is the evidence for reliability and validity of behavior evaluations for shelter dogs? a prequel to “no better than flipping a coin”. *Journal of Veterinary Behavior* (2019).
- [17] Borsboom, D., Mellenbergh, G. J. & Van Heerden, J. The concept of validity. *Psychological review* **111**, 1061 (2004).
- 610 [18] Borsboom, D., Cramer, A. O., Kievit, R. A., Scholten, A. Z. & Franić, S. The end of construct validity. In *The Concept of Validity: Revisions, New Directions and Applications, Oct, 2008* (IAP Information Age Publishing, 2009).
- [19] Maul, A., Irribarra, D. T. & Wilson, M. On the philosophical foundations of psychological measurement. *Measurement* **79**, 311–320 (2016).
- 615 [20] Patronek, G. J. & Bradley, J. No better than flipping a coin: Reconsidering canine behavior evaluations in animal shelters. *Journal of Veterinary Behavior* **15**, 66–77 (2016).
- [21] Valsecchi, P., Barnard, S., Stefanini, C. & Normando, S. Temperament test for re-homed dogs validated through direct behavioral observation in shelter and home environment.  
620 *Journal of Veterinary Behavior* **6**, 161–177 (2011).
- [22] ASPCA. Position statement on shelter dog behavior assessments. <https://www.aspc.org/about-us/aspc-policy-and-position-statements/position-statement-shelter-dog-behavior-assessments> (2018). Online; accessed 13-May-2020.
- 625 [23] Clay, L. *et al.* In defense of canine behavioral assessments in shelters: Outlining their positive applications. *Journal of Veterinary Behavior* (2020).

- [24] Goold, C. M. & Newberry, R. C. Aggressiveness as a latent personality trait of domestic dogs: testing measurement invariance and local independence. *PLoS One* (2017).
- [25] Goold, C. & Newberry, R. C. Modelling personality, plasticity and predictability in shelter dogs. *Royal Society open science* **4**, 170618 (2017).
- [26] Dingemanse, N. J., Kazem, A. J., Réale, D. & Wright, J. Behavioural reaction norms: animal personality meets individual plasticity. *Trends in ecology & evolution* **25**, 81–89 (2010).
- [27] Cleasby, I. R., Nakagawa, S. & Schielzeth, H. Quantifying the predictability of behaviour: statistical approaches for the study of between-individual variation in the within-individual variance. *Methods in Ecology and Evolution* **6**, 27–37 (2015).
- [28] Vonesh, E. F., Greene, T. & Schluchter, M. D. Shared parameter models for the joint analysis of longitudinal data and event times. *Statistics in medicine* **25**, 143–163 (2006).
- [29] Tseng, C.-h., Elashoff, R., Li, N. & Li, G. Longitudinal data analysis with non-ignorable missing data. *Statistical methods in medical research* **25**, 205–220 (2016).
- [30] Lambert, D. Zero-inflated poisson regression, with an application to defects in manufacturing. *Technometrics* **34**, 1–14 (1992).
- [31] Kelley, M. E. & Anderson, S. J. Zero inflation in ordinal data: incorporating susceptibility to response through the use of a mixture model. *Statistics in medicine* **27**, 3674–3688 (2008).
- [32] Bagozzi, B. E. & Mukherjee, B. A mixture model for middle category inflation in ordered survey responses. *Political Analysis* **20**, 369–386 (2012).
- [33] Kubinec, R. Generalized ideal point models for time-varying and missing-data inference. *OSF Preprint* (2019).
- [34] Kruschke, J. *Doing Bayesian data analysis: A tutorial with R, JAGS, and Stan* (Academic Press, 2014).
- [35] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria (2019).
- [36] Stan Development Team. Stan modeling language user’s guide and reference manual, version 2.18.0 (2018).
- [37] Stan Development Team *et al.* Cmdstan: the command-line interface to stan. version 2.23.0. (2018).



- [38] IntHout, J., Ioannidis, J. P., Rovers, M. M. & Goeman, J. J. Plea for routinely presenting prediction intervals in meta-analysis. *BMJ open* **6**, e010247 (2016).
- 660 [39] Wang, C.-C. & Lee, W.-C. A simple method to estimate prediction intervals and predictive distributions: Summarizing meta-analyses beyond means and confidence intervals. *Research synthesis methods* **10**, 255–266 (2019).
- [40] Gates, M. C., Zito, S., Thomas, J. & Dale, A. Post-adoption problem behaviours in adolescent and adult dogs rehomed through a new zealand animal shelter. *Animals* **8**,  
665 93 (2018).
- [41] Gelman, A. & Hill, J. *Data analysis using regression and multilevelhierarchical models*, vol. 1 (Cambridge University Press New York, NY, USA, 2007).
- [42] McElreath, R. & Koster, J. Using multilevel models to estimate variation in foraging returns. *Human nature* **25**, 100–120 (2014).
- 670 [43] Kerr, N. L. Harking: Hypothesizing after the results are known. *Personality and social psychology review* **2**, 196–217 (1998).
- [44] McElreath, R. *Statistical rethinking: A Bayesian course with examples in R and Stan* (CRC press, 2020).