

1 Counter culture: Causes, extent and 2 solutions of systematic bias in the analysis 3 of behavioural counts

4 Joel L. Pick^{1,2,3}, Nyil Khwaja^{2,4}, Michael A. Spence^{2,5}, Malika Ihle^{2,6}, and
5 Shinichi Nakagawa¹

6 ¹School of Biological, Earth and Environmental Sciences, University of New South Wales,
7 Randwick NSW 2052, Sydney, Australia

8 ²Department of Animal and Plant Sciences, University of Sheffield, Western Bank,
9 Sheffield S10 2TN, UK

10 ³Institute of Ecology and Evolution, University of Edinburgh, Charlotte Auerbach Road,
11 Edinburgh, EH9 3FL, UK

12 ⁴School of Biological Sciences, University of Canterbury, Christchurch, New Zealand

13 ⁵Centre for Environmental Fisheries and Aquaculture Science, Lowestoft, Suffolk, NR33
14 OHT, UK

15 ⁶Department of Psychology, Ludwig-Maximilians-Universität München, Leopoldstr. 13,
16 80802, München, Germany

17 Corresponding author:

18 Joel L. Pick¹

19 Email address: joel.l.pick@gmail.com

20 ABSTRACT

21 We often quantify the rate at which a behaviour occurs by counting the number of times it occurs within
22 a specific, short observation period. Measuring behaviour in such a way is typically unavoidable but
23 induces error. This error acts to systematically reduce effect sizes, including metrics of particular interest
24 to behavioural and evolutionary ecologists such as R^2 , repeatability (intra-class correlation, ICC) and
25 heritability. Through introducing a null model, the Poisson process, for describing the frequency of
26 behaviour, we give a mechanistic explanation of how this problem arises and demonstrate how it makes
27 comparisons between studies and species problematic, because the magnitude of the error depends on
28 how frequently the behaviour has been observed as well as how biologically variable the behaviour is.
29 Importantly, the degree of error is predictable and so can be corrected for. Using the example of parental
30 provisioning rate in birds, we assess the applicability of our null model for describing the frequency of
31 behaviour. We then survey recent literature and demonstrate that the error is rarely accounted for in
32 current analyses. We highlight the problems that arise from this and provide solutions. We further discuss
33 the biological implications of deviations from our null model, and highlight the new avenues of research
34 that they may provide. Adopting our recommendations into analyses of behavioural counts will improve
35 the accuracy of estimated effect sizes and allow meaningful comparisons to be made between studies.

36 INTRODUCTION

37 Behaviour is frequently quantified by counting the number of times a specific event occurs within
38 an observation period. This includes studies of parental care (e.g. number of feeding visits), social
39 interactions such as allopreening/grooming or aggression (e.g. number of interactions) and mate choice
40 (e.g. number of copulations, courtship behaviours). Typically, when behaviour is quantified in such a
41 way, we do not observe the total time over which a behaviour takes place, and thus, all occurrences of the
42 behaviour. Rather we sample a shorter time period in order to calculate a representative 'rate' at which the
43 behaviour occurs. Take the example of parental provisioning behaviour in birds - although the nestling
44 period may last 2 or more weeks, researchers typically record the feeding visits that occur in a shorter
45 period of time, often in a 1 or 2 hour period (Murphy et al., 2015). Measuring behaviour in this way

46 is often unavoidable for practical reasons, nevertheless we accept that there will be some error in the
47 quantification of a behaviour. Intuitively we also know that the longer we observe a behaviour, the better
48 the representation of that behaviour we will have (i.e. diminishing error with increasing sampling effort;
49 see Murphy et al., 2015; Lendvai et al., 2015; Sánchez-Tójar et al., 2018, for empirical evidence of this).

50 Evolutionary and behavioural ecologists are often interested in quantifying the total amount of
51 variation that is due to differences between individuals and environments. We therefore frequently use
52 metrics such as R^2 (proportion of total variance explained by a particular model), repeatability (the
53 proportion of total variance due to individual identity effects; also known as the intraclass correlation
54 coefficient - ICC) or heritability (proportion of total variance due to additive genetic effects) or, more
55 broadly, quantify standardised effect sizes (Nakagawa and Cuthill, 2007; Nakagawa and Schielzeth, 2010,
56 2013). However, commonly used methods for analysing variation in behavioural counts fail to distinguish
57 between the contribution of biological variation and the error introduced by sampling. Without taking this
58 into account, we will both systematically *underestimate* effect sizes, and limit our ability to compare such
59 metrics between studies, due to variation in sampling effort.

60 Here, we outline a broad model for thinking about how behavioural count data arise, which demon-
61 strates the inevitability of systematic error in such data. We suggest a simple null model for describing
62 these processes and discuss the problems associated with ignoring the stochasticity inherent in this data.
63 Using the example of provisioning rate in birds, a frequently used, count-based measure of parental care,
64 we show how behavioural counts fit our null model well. Through a literature survey we show how
65 widespread the problem of not accounting for this stochastic error is. Finally we present some solutions to
66 this problem. Despite the large focus of this article on provisioning, all the theoretical and practical issues
67 and solutions described here are directly relevant to other behaviours sampled by counting events in a
68 restricted time window.

69 BEHAVIOUR AS A POINT PROCESS

70 Let us take a general example to describe our behavioural count data. Imagine that we want to describe
71 the factors affecting the frequency (or rate) at which bees enter a nest. Here each arrival represents the
72 behaviour of interest occurring. To do this, we watch the entrance to a nest for set observation periods (e.g.
73 1 hour). There will be many factors influencing the length of the intervals between arrivals (interval length),
74 only a few of which we will realistically be able to quantify (e.g. biotic factors such as the abundance and
75 distribution of food/competitors/predators and abiotic factors such as temperature/rainfall/wind/humidity
76 at each moment during the this observation). As we are unable to describe the complex processes that
77 lead to the timing of this behaviour, we can therefore describe the arrival of these bees (or more generally
78 the occurrence of behaviour), as a stochastic process through time. Note that we are not suggesting that
79 the behaviour arises through a stochastic process, rather that we can *describe* it as a stochastic process;
80 although the behaviour may be completely deterministic, we do not have the information needed to
81 describe it in this way. We therefore need a model that describes the stochastic distribution of events
82 through time. This is broadly known as a point process, a probabilistic model for describing points (or
83 events) in some space, for example the distribution of points occurring on a straight line or analogously as
84 events occurring through time (Figure 1, and see Table 1 for glossary of terms).

85 Introducing the Poisson process

86 A commonly used point process is the Poisson process (Daley and Vere-Jones, 2003). It has a single
87 parameter, the arrival rate (λ) for a given unit of time (e.g. 10 arrivals/hour). By using counts of behaviour
88 (from a sampling period) as a representative measure of that behaviour, we are making the assumption that
89 there is an underlying rate, which we try to capture during the observation. As we want to understand the
90 factors affecting the rate at which bees enter a nest, it is this parameter, λ , that we are broadly interested
91 in estimating (note that we never observe this rate directly, it is an underlying, latent property of this
92 process).

93 The simplest Poisson process is the homogeneous Poisson process, which assumes that the arrival
94 rate (λ) is constant through time. Given the simplicity of this process, we believe that it is a highly
95 suitable starting point or null model with which to describe behavioural counts. It also has several useful
96 properties, which we can compare with real data to assess the suitability of this model (explored in *Does a*
97 *Poisson process provide a good description of behaviour?*). First, it assumes that the probability of arrival
98 is constant over time, which results in the interval lengths within an observation being exponentially

99 distributed. Second, a Poisson process results in a predictable amount of variation in the number of visits
 100 observed (per observation) across multiple observations; these counts (for each observation) follow a
 101 Poisson distribution, with mean (and variance) equal to the length of the observation period (t) multiplied
 102 by the arrival rate, or λt . We will describe this mean as the *expected* number of arrivals for an observation
 103 (c.f. de Villemereuil et al., 2016). In other words, for every observation there is an underlying arrival
 104 rate leading to an expected number of arrivals, but we observe this with a certain degree of (quantifiable)
 105 error caused by the stochastic nature of the process (Figure 2A). For example, if λ is 10 arrivals/hour,
 106 then the expected number of arrivals in one hour would be 10, two hours 20 etc, but we would observe
 107 considerable variation around these expectations. Formulaically we describe this as $y \sim \text{Poisson}(\lambda t)$,
 108 where y is the observed counts. We will refer to this error as *stochastic* error. In the statistical literature,
 109 this would be referred to as aleatory uncertainty (Kiureghian and Ditlevsen, 2009). This error is analogous
 110 to measurement error, with similar consequences (outlined below).

111 Now imagine that we want to compare different bee nests (for instance we may be interested in the
 112 differences in arrival rate due to factors such as the nest's size, distance to food etc.). If all nests had
 113 the same arrival rate, then by watching them for the same observation period, the observed number of
 114 arrivals across these observations would be Poisson distributed (Figure 2A). In other words, we would
 115 still observe variation in the number of arrivals between observations of different nests, even if there is no
 116 variation in their arrival rates (Figure 2A). If the arrival rates differ between nests (i.e. variation in the
 117 expected number of arrivals, due to factors such as nest size etc) the combined observations would be
 118 over-dispersed (i.e. more variation present than explained simply by the Poisson distribution; Figure 2B).

119 **Diminishing stochastic error**

120 As mentioned above, it seems intuitive that longer observations result in lower error (Lendvai et al., 2015;
 121 Morvai et al., 2016). As both the mean number of observed arrivals across all observations increases (for
 122 example through longer observations) and with greater variation in arrival rates, the amount of stochastic
 123 error, relative to total observed variance, diminishes. Let's look more carefully at why this is the case.

124 First, following the law of total variance, the total observed variance in the number of arrivals across
 125 observations (σ_y^2) is equal to the expectation of the stochastic variance (σ_{stoc}^2) plus the variance in expected
 126 number of arrivals ($\sigma_{\lambda t}^2$)

$$\sigma_y^2 = \sigma_{stoc}^2 + \sigma_{\lambda t}^2 \quad (1)$$

127 Second, the stochastic variance and variance in the expected number of arrivals do not scale in the same
 128 way. This is most easily demonstrated by thinking of variability in terms of the coefficient of variation, a
 129 dimensionless measure of variability ($CV_x = \sigma_x/\bar{x}$, where \bar{x} and σ_x are the mean and standard deviation,
 130 respectively, of x). The variation in expected number of arrivals, when measured as CV, remains constant
 131 as the mean number of arrivals increases. This is because the standard deviation scales directly with
 132 the mean (e.g. $\sigma_{2x} = 2\sigma_x$). So, if the observation period was twice as long, all the expected number of
 133 visits across observations would double, as would their SD, whilst the CV remains the same (Figure 3A;
 134 $2\sigma_x/2\bar{x} = \sigma_x/\bar{x}$). Put as an equation, if the observation period (t) is constant across observations,

$$CV_{exp} = \sigma_{\lambda t}/\bar{\lambda}t = t\sigma_{\lambda}/t\bar{\lambda} = \sigma_{\lambda}/\bar{\lambda} \quad (2)$$

135 and so changing the observation period does not change the CV in the expected number of visits. In
 136 contrast, due to the nature of the Poisson distribution (i.e. the mean equals the variance), as the mean
 137 (number of arrivals) increases, the Poisson distributed stochastic error becomes relatively less variable (i.e.
 138 CV decreases; $CV_{stoc} = \sqrt{\bar{\lambda}t}/\bar{\lambda}t$; Figure 3A). In other words, the stochastic error diminishes relative to the
 139 variation in expected number of arrivals as the sampled number of arrivals becomes larger.

140 Given that $\sigma_{stoc}^2 = \bar{\lambda}t$ and that $\sigma_{\lambda t}^2 = (\bar{\lambda}t CV_{exp})^2$ we can rewrite equation 1 as:

$$\sigma_y^2 = \bar{\lambda}t + (\bar{\lambda}t CV_{exp})^2 \quad (3)$$

141 We can now see that the variation due to stochastic error is equal to the mean, whilst the variation in the
 142 expected number of visits across observations is a function of both the expected CV (a constant) and
 143 the square of the mean. Therefore, as the mean number of observed visits increases, the variance due
 144 to the expected number of arrivals increases exponentially compared with stochastic error. This results

145 in a rapid decrease in the observed CV and an increase in the proportion of the observed variation due
146 to variation in the expected number of arrivals (Figures 3B, C). In other words, the amount of times a
147 behaviour occurs in an observation period directly affects the confidence we can have in quantifying
148 the rate at which it occurs. Equally, the greater the amount of underlying variation in arrival rates (i.e.
149 expected CV; shown by different lines in Figures 3B, C), the lower the impact the stochastic error has and
150 the more the observed variation reflects this expected variation in arrival rate (Figure 3C).

151 This stochastic error can be seen as analogous to measurement error; if we take a variable measured
152 with a large amount of error, averaging over an increasing number of measurements would give a better
153 estimate. Similarly, the more events we observe, the more the stochastic error is averaged over, and the
154 greater precision we have in our estimate of λ .

155 **Dynamic rates**

156 A homogeneous Poisson process assumes that the rate at which a behaviour occurs is the same throughout
157 the observation. It is, however, possible that individuals adjust their behaviour at a very fine scale (e.g.
158 Johnstone et al., 2014; Schlicht et al., 2016), meaning that the expected rate changes during the observation.
159 A changing arrival rate does not, however, violate the assumptions of a general Poisson process. Indeed,
160 there exist models that allow the rate to be dynamic (e.g. inhomogenous Poisson processes Heuer et al.,
161 2010). For example, in a Cox point processes, a generalisation of the Poisson process, the arrival rate
162 is also assumed to be stochastic, and modelled as a latent variable (Blackwell et al. 2016; Spence et al.
163 2021; see also Johnstone et al. 2014 for a related example in behavioural ecology). When not modelled,
164 such dynamically changing behaviour will add further error to the estimation of the overall rate at which
165 a behaviour occurs. Consequently, if the behaviour of interest is extremely dynamic, we would not
166 expect to find anything we measure on a broad scale to correlate with it. However, many behaviours
167 have been found to consistently differ between individuals, as shown in the recent flurry of studies on
168 animal personality (Bell et al., 2009; Beekman and Jordan, 2017), and are frequently found to vary among
169 different environments. Such findings indicate that there is a consistent rate that can be measured over
170 these time periods in many behaviours.

171 **Refractory periods and the non-independence of visits**

172 A homogeneous Poisson process also assumes that, at any time point, the time to next arrival is independent
173 of when the previous arrival occurred (i.e. the likelihood of an event occurring is constant over time).
174 This is known as the Markov property, and results in an exponential distribution of interval lengths (with
175 a modal interval length of 0). Up until now we have been considering arrivals of bees at a nest, which are
176 likely to be largely independent of each other. When considering the behaviour of a single individual,
177 however, this may seem unrealistic. It is important to note that when considering this assumption of
178 independence, it does not matter if we are watching an individual, a pair, a group or a set of unique
179 individuals. The assumption of independence relates simply to whether the timing of one event affects the
180 occurrence of the next event, and not to the individual that does the event (although the chance that the
181 assumption of independence is violated may increase with fewer individuals).

182 Let us consider then, the arrival rate to the nest of an individual bee (rather than the overall arrival
183 rate at the nest). As bees have to find food and return to the nest, the probability of this individual bee
184 arriving at the nest is likely to be lower just after its last arrival occurred, meaning the probability of the
185 bee arriving changes over time. This is known as a refractory period (i.e. a period in which a behaviour is
186 unlikely to reoccur). A refractory period can be described in a point process with an additional parameter.
187 As stated above, a Poisson process assumes that the interval lengths follow an exponential distribution,
188 with a modal interval length of 0. The exponential distribution is a special case of the gamma distribution,
189 in which one of its two parameters, α , is fixed to 1 (no refractory period). When $\alpha > 1$, we have a point
190 process with a refractory period. α describes the refractory period in terms of the expected interval length.
191 The refractory period itself can be more intuitively thought of as the mode of the gamma distribution (see
192 Supplementary Material S1).

193 As the interval lengths no longer follow an exponential distribution when there is a refractory period,
194 the number of arrivals in a given period of time would no longer follow a Poisson distribution, violating
195 the assumptions of the Poisson process. A refractory period reduces the amount of stochastic error for a
196 particular mean (i.e. underdispersion with respect to a Poisson distribution) in a predictable way, and so
197 can be modelled using the additional α parameter to describe the relative extent of the refractory period
198 (see Supplementary Material S1). Although a refractory period results in less stochastic error for given

199 mean number of observations, it is important to note that the stochastic error will always be proportional
200 to this mean, and so the relative amount of stochastic error will still diminish as more events are observed,
201 or with more biological variation. In the context of the equations above, with a refractory period σ_{stoc}^2
202 becomes $\bar{\lambda}t/\alpha$ rather than $\bar{\lambda}t$ ($1/\alpha$ is commonly referred to as ϕ or the dispersion parameter in some
203 analytical models). Furthermore, small refractory periods do not lead to substantial deviations away from
204 a Poisson process.

205 In behavioural data, there are several instances where such refractory periods may exist, for example,
206 if an individual has to do something before the behaviour reoccurs. As in the example above, this might
207 be seen in arrival or feeding rates, especially if foraging sites are located far away from the nests. Sexual
208 behaviours are also likely to show such refractory periods; this is well studied in rats and humans for
209 example (Levin, 2009), and it is likely that there is a minimum interval length between copulations in
210 many other species. Quantification of refractory periods is therefore important, but to our knowledge
211 has not yet been systematically investigated for any behavioural trait measured this way (at least in the
212 context of behavioural ecology). The little quantitative information about the presence, and length, of
213 refractory periods makes it very difficult to judge their impact, at least currently.

214 **Embracing a null model**

215 *... there is no need to ask the question "Is the model true". If "truth" is to be the "whole*
216 *truth" the answer must be "No". The only question of interest is "Is the model illuminating*
217 *and useful?"*

218 George Box, 1979

219 Here we argue that a Poisson process is a good *null* model to describe the stochastic nature of
220 behaviours sampled in specific periods of time (a type of point process) for the following 3 reasons.
221 First, it is simple and tractable, and it allows us to account for how stochastic error predictably changes
222 with sampling effort. We acknowledge that it may not be the perfect model for such behaviours in all
223 circumstances (as discussed above). However, this is not the purpose of a null model. Currently, we have
224 no clear null model, as generally we simply ignore the presence of this stochastic error (see literature
225 survey below), which is induced by the processes underlying behavioural count data. As we will discuss
226 below, ignoring this stochastic error (the presence of which we believe to be undeniable, the form/extent
227 of which can be the subject of debate) leads to systematic error in the analysis of behavioural counts.
228 Second, deviation from this model gives us valuable information. We believe that assuming (and more
229 importantly understanding) such a null model gives us insight into the processes underlying the data,
230 whether deviations from this model occur, and if so what they may represent. Many fields have embraced
231 null models (for example, the ideal-free distribution in behavioural ecology or the Hardy-Weinberg
232 Equilibrium in population genetics), and it is standard practice to quantify deviation from these models.
233 Finally, such a null model would force us to confront what assumptions we are making when analysing
234 our behavioural count data.

235 **PROBLEMS WITH IGNORING STOCHASTIC ERROR**

236 Generally, stochastic error induces the same analytical problems as measurement error. Although these
237 have been covered elsewhere (e.g. Freckleton, 2011; Garamszegi, 2016; Ponzi et al., 2018; Dingemanse
238 et al., 2021), they have not been discussed in the context of behavioural count data, and so we briefly
239 outline the problems here for completeness.

240 **Analysing variation in behavioural counts**

241 When stochastic error is not explicitly modelled, it is included in the residual, unexplained variation
242 (Figure 4). This imposes an upper limit to the variance in, for example, arrival rate that can be explained
243 (Figure 3C), because the stochastic error will always remain unexplained, unless explicitly accounted for
244 (Figure 4; Nakagawa and Schielzeth, 2013; Nakagawa et al., 2017a). To demonstrate this, let us return to
245 our population of bee nests where the mean (\pm SD) arrival rate across nests is 10 ± 3 arrivals/hour (i.e. an
246 expected CV of 0.3), and 50% of this expected variation in arrival rate is due to consistent differences
247 between nests (i.e. repeatability (ICC) on the expected scale = 0.5). Assuming that arrival rates in this
248 population are well described by a Poisson process, a study that observed nests for 2 hours (Study A;
249 Figure 4), would observe an average of 20 ± 7.5 arrivals per observation and estimate a repeatability of

250 0.32, just over half of the actual repeatability. Not accounting for this stochastic error leads, therefore, to
251 a general underestimation of underlying effect sizes.

252 Different studies will also vary in the mean number of observed arrivals and/or the underlying variation
253 in expected arrival rates, through having different observation periods, or simply because of intrinsic
254 differences among populations. As the proportion of total variance due to stochastic error is dependent on
255 both of these factors (Figure 3C), metrics that relies on the estimation of total variance (e.g. standardised
256 effect sizes, ICC and R^2) are not comparable between studies, when not accounting for this changing
257 stochastic error. Imagine that two more studies (Studies B and C) are performed in the population
258 described above, but with shorter observation periods (60 and 30 mins, respectively), averaging 10 and
259 5 arrivals per observation. As a different amount of stochastic error was observed in both cases, the
260 resulting repeatabilities would be much lower still, 0.24 and 0.16 respectively (Figure 4). Effect sizes,
261 therefore, may differ between studies due to both the intrinsic characteristics of the population and the
262 sampling effort. Note that these calculations assume of an underlying Poisson process. If, for example,
263 the refractory period was to differ between different studies then, without accounting the differing form of
264 stochastic error, the results would also not be comparable.

265 **The low predictive power of behaviour**

266 Behavioural count data typically correlates poorly with other variables. However, because of the potentially
267 large proportion of observed variation that is due to stochastic error, the observed number of events is
268 constrained in how much variation in another trait it can explain (Figures 3C, 4). For example, arrival rate
269 estimated from a short observation period will correlate poorly with the underlying arrival rate, due to this
270 stochastic error (Lendvai et al., 2015; Morvai et al., 2016). Consequently, this measure may explain little
271 variation in another variable - even if it actually has had a strong effect. Moreover, stochastic error in
272 one predictor variable can have a large effect on the parameter estimates of other covariates in the model,
273 as the covariance between different parameters is not properly estimated, creating potentially spurious
274 relationships between predictor variables and the response variable (Freckleton, 2011). Note that these
275 effects are a general consequence of any kind of measurement error in predictor variables, but will be
276 particularly pronounced with the level of error seen in count data.

277 **DOES A POISSON PROCESS PROVIDE A GOOD DESCRIPTION OF BE-** 278 **HAVIOUR?**

279 Up until now we have been focusing on a general example whilst arguing that a Poisson process represents
280 a suitable descriptive model for behavioural count data. Now we will take a frequently used behavioural
281 count - provisioning rate - to demonstrate the utility of this model and highlight the extent of the problems
282 caused by not taking such stochastic error into account. Parental provisioning rate (measured as the
283 number of feeding visits within a certain unit of time) is often used as a quantitative assessment of
284 parental investment in birds and analyses of provisioning rate have contributed a considerable amount
285 to our understanding of parental care (e.g. Harrison et al., 2009). As a Poisson process makes certain
286 assumptions, we can compare the patterns we see in observed data with those we expect from a Poisson
287 process, to assess how good a model this is for describing our data. There are 3 important patterns which
288 can emerge.

289 First, a Poisson process assumes that visits to the nest are independent from each other, in others words
290 there is no refractory period. Note that the assumption of independence is not violated by our watching
291 the same individual, rather referring to the probability of a visit occurring depending on when the last
292 visit occurred. Although this is perhaps the most obvious way in which provisioning rate could violate
293 the assumptions of a Poisson process, at the moment, there is little evidence that substantial refractory
294 periods exist for provisioning rate. For example, distributions of inter-visit interval lengths that appear
295 close to exponential (as expected from a Poisson process) have been observed in several species (Great
296 tit (*Parus major*) - Johnstone et al. 2014; Acorn Woodpecker (*Melanerpes formicivorus*) - Figure 3 in
297 Koenig and Walters 2016; Red Winged Blackbird (*Agelaius phoeniceus*) - Figure 2 in Westneat et al. 2013;
298 Pied Flycatcher (*Ficedula hypoleuca*) - Figure S1 in Westneat et al. 2017; Chestnut-crowned Babbler
299 (*Pomatostomus ruficeps*) - Savage et al. 2017; House Sparrow (*Passer domesticus*) - Figure S1 in Ihle
300 et al. 2019). Furthermore, many studies analyse per nest visit rates (i.e. with two or more parents/carers),
301 in which case refractory periods are likely to be extremely low.

302 Second, from a Poisson process we would expect Poisson distributed error. With additional factors
303 influencing provisioning rate (due to individual, brood or environmental characteristics), we would
304 see additional variation (i.e. overdispersion with respect to a Poisson distribution; Figure 2) and so
305 the variance in the number of observed visits between observations would be greater than the mean.
306 Observing overdispersion does not necessarily mean that the stochastic error is Poisson distributed; we
307 could still observe more variation than expected in a Poisson distribution, if for example the stochastic
308 error was lower (due to a refractory period) and biological variation greater. The variance being more than
309 or equal to the mean is rather a minimum requirement for Poisson distributed error to exist. Consistent
310 underdispersion with respect to a Poisson distribution across population level estimates of variation, on
311 the other hand, would indicate that the error was not Poisson distributed, violating the assumptions of our
312 null model. As shown below (see Literature Survey), overdispersion is consistently found across studies
313 and species in provisioning rate, which is consistent with our model of a Poisson process and additional
314 between-observation variation in expected provisioning rates.

315 Finally, we would expect to see a dramatic decrease in the relative variability of provisioning data with
316 an increase in the mean number of observed visits (Figure 3A). Using recently published data (Lendvai
317 et al., 2015), we can see this systematic decrease in CV with increasing observation time (Figure 3D; see
318 Supplementary Material S2 for further details), in accordance with theoretical predictions. Concurrently
319 we would also expect that metrics such as repeatability would increase with the mean number of observed
320 visits (Figure 3C). Again, we can show this effect empirically using the data of Lendvai et al. (2015). By
321 calculating the repeatability of provisioning rate using the same overall total time period, but split into
322 differently sized observations periods (and not correcting for this stochastic error), we indeed find that
323 repeatability dramatically increases with observation period (i.e. with mean number of observed visits;
324 Figure 3E, Supplementary Material S2), in line with expectations from a Poisson process.

325 Together this evidence demonstrates that provisioning data has a high level of stochastic error, the
326 magnitude of which is in line with that predicted by a Poisson process. Therefore, a Poisson process
327 seems a highly suitable model for provisioning rate. We should stress, however, that the validity of these
328 assumptions should be assessed on a behaviour and study specific basis.

329 **ASSESSING THE SIZE OF THE PROBLEM - LITERATURE SURVEY**

330 Recent work has suggested that only a small amount of variation in provisioning rate is generally explained
331 by individual, brood or environmental characteristics (Williams, 2012; Williams and a. Fowler, 2015), a
332 trend that is commonly found across behavioural traits (e.g. low repeatabilities; Bell et al., 2009; Wolak
333 et al., 2012). It has also been suggested that provisioning rate often has little or no detectable effects on
334 offspring phenotype (e.g. fledgling size, survival etc; Schwagmeyer and Mock, 2008; Williams, 2012;
335 Williams and a. Fowler, 2015), bringing into question its utility as an indicator of parental investment.
336 However, such conclusions may arise from failing to account for the presence of stochastic error. How
337 much of a problem these issues present depends largely on both the sampling effort employed in such
338 studies, how variable this effort is among studies, and whether or not the presence of stochastic error is
339 accounted for. In order to ascertain the breadth of the problems outlined above, we conducted a survey of
340 papers analysing provisioning rate (measured as the number of visits), published in 2015/16. We did not
341 intend the search to be exhaustive, rather to generate a representative set of recent papers on provisioning
342 rate.

343 **Survey Methods**

344 On 13/12/2016 JLP searched Web of Science using the search term:

```
345 (TS=("visit rate" OR "number of visit*" OR "nest visit*" OR "provisioning"  
346 OR "feeding rate" OR "parental care" OR "number of feed*") AND TS=(chick*  
347 OR nest* OR fledg* OR offspring) AND TS=(*bird* OR passerine* OR avian  
348 OR chick*) NOT TS=(veterinary OR chicken* OR hen* OR broiler* OR poult*  
349 OR layer* OR "Japanese quail*" OR turkey* OR chickpea* OR pollin*)) AND  
350 PY=(2015 OR 2016) AND LANGUAGE: (English) AND DOCUMENT TYPES: (Article)
```

351 returning 289 papers. JLP and MI screened the abstracts to exclude any papers that did not relate to
352 provisioning (or show primary data e.g. reviews). JLP and NK read the remaining 143 papers, looking
353 specifically for studies that measured provisioning as number of visits, as opposed to inter-visit intervals

354 or quantity of food. Although we did not target these studies, we included any studies that measured
 355 incubation feeding by males. This process returned 81 studies. At a reviewers request, JLP repeated this
 356 process for papers published in 2022, searching Web of Science on 05/12/2022, returning 147 papers,
 357 which reduced to 46 after abstracts screening, 29 of which analysed provisioning data.

358 From the 81 papers from 2015/16, JLP and NK extracted the study species, the method of data
 359 collection (direct, video, RFID etc) and length of observation period (hours). We also extracted summary
 360 data for each analysis that was conducted using provisioning rate, totalling 427 analyses across all
 361 the studies. For each analysis, we recorded whether provisioning was used as a response or predictor
 362 variable. If analysed as a response, we recorded what error distribution was used e.g. Gaussian (including
 363 linear regressions, t-tests, ANOVAs, correlations), Poisson, Negative Binomial, non-parametric etc. If
 364 not otherwise stated, we assumed Gaussian error distribution was used, because this is the default in
 365 most statistical software. We recorded whether the number of visits itself was analysed, or whether
 366 it was transformed into a rate (e.g visits per hour or visits per chick). The latter of these metrics
 367 (response/predictor, distribution, number/rate) were also extracted from the 29 papers from 2022 by JLP,
 368 giving 79 analyses.

369 For each 2015/16 analysis, we then extracted, where possible, mean and SD of provisioning data used
 370 in that analysis (note some analyses used the same data, so we recorded how many times that data set was
 371 used and in what ways). If SE and sample size (N) were presented SD was calculated as $SD = SE\sqrt{N}$. If
 372 range and/or inter-quartile range were presented then a SD was derived following the formulas presented
 373 in Wan et al. (2014). If not presented in the text, mean and SD or SE were extracted from figures using
 374 the *digitize* and *metaDigitise* packages in *R* (Poisot, 2011; Pick et al., 2019), or from raw data if available.
 375 If parameter estimates from models were presented, the intercept was taken as a measure of the mean
 376 provisioning rate (and back-transformed if necessary), as long as any covariates were centered, but the
 377 associated SE was not used. We aimed to collect means and SDs that were representative of the data that
 378 was analysed. This meant that when means and SDs from different groups included in the same analyses
 379 were presented (for example, means presented per sex, but data analysed across both sexes), we pooled
 380 these M estimates (from samples x_1 through x_M) according to the formulas:

$$\bar{x} = \frac{\sum_i N_{x_i} \bar{x}_i}{\sum_i N_{x_i}} \quad (4)$$

$$\sigma_x = \sqrt{\frac{1}{\sum_i N_{x_i} - 1} \left(\sum_i [(N_{x_i} - 1)\sigma_{x_i}^2 + N_{x_i}\bar{x}_i^2] - \left[\sum_i N_{x_i} \right] \bar{x}^2 \right)}. \quad (5)$$

381 Where the mean (\pm SD) provisioning was presented as a rate (i.e. had been transformed from the scale on
 382 which it was collected), we back-transformed these to their original scale, e.g. if the observation period
 383 was 4 hours, but provisioning was presented as visits/hour, we multiplied both the mean and SD by 4. For
 384 estimates where the observation period varied, we corrected them relative to the mean observation period
 385 (or mode if this was presented). Similarly, when rates were presented as per chick, and the mean number
 386 of chicks was presented, then we back calculated the mean, but not SD. For estimates for which we had
 387 SD, and were not presented as per chick, we calculated the expected CV as

$$CV_{exp} = \frac{\sqrt{\sigma_x^2 - \bar{x}}}{\bar{x}} \quad (6)$$

388 . For several estimates, especially when the mean is low the expected CV cannot be calculated, and the
 389 mean is larger than the variance (see Supplementary Material S5). To these estimates we gave a value of
 390 0. We then calculated the proportion of observed variation due variation in expected provisioning rate as

$$\frac{\bar{x}CV_{exp}^2}{1 + \bar{x}CV_{exp}^2} \quad (7)$$

391 The papers included in the study, and the data extracted from them, are presented in Table S1.

392 Survey Results

393 Our survey of papers published in 2015/16 found 81 studies of 64 species that fitted our criteria, containing
 394 427 analyses of provisioning rate (343 as a response and 84 as a predictor) and our 2022 survey found

395 29 studies containing 79 analyses of provisioning rate (50 as a response and 29 as a predictor). For 350
396 of the 2015/16 analyses we could extract the mean number of observed visits for the data used in the
397 analysis, which was typically low (median = 8.41, Figure 5). For 301 analyses we were also able to
398 extract a standard deviation of provisioning rate, and so calculate expected CV. For 34 of these analyses,
399 it was not possible to calculate expected CV, and so it was assumed to be 0 (Supplementary Material S5).
400 Expected CV ranged from 0 to 1.234 (median=0.449).

401 Using the expected CV and mean number of observed visits for each sample, we were able to calculate
402 the proportion of observed variation in provisioning rate that is due to expected, biological, variation (see
403 Supplementary Material S3). Across these estimates, the median proportion was 0.627 (Figure 5). This
404 represents the maximum effect size (e.g. R^2 or ICC) that can be estimated from this data, if sampling
405 error is not accounted for (i.e. if the true ICC was 1 the estimated ICC would be 0.627, and if the true
406 ICC were 0.5, then the estimated ICC would be 0.3135 etc.). Of the 427 analyses, only 7% and 22% (in
407 2015/16 and 2022, respectively) modelled the stochastic error by assuming a Poisson error distribution
408 when provisioning rate was analysed as a response variable (see below, note no other corrections were
409 used), whilst no study in either survey accounted for stochastic error when modelling provisioning rate
410 as a predictor variable. Effect sizes from these studies are therefore consistently, and often substantially,
411 underestimated. Repeatabilities, for example, would be underestimated on average by 37%, as would
412 the effect of provisioning rate on other variables such as chick mass or survival. Moreover, because of
413 the large variation in this proportion among studies (range: 0 - 0.982), interpretation of, and comparison
414 amongst, effect size estimates from these studies is not meaningful, as any differences may simply be due
415 to methodology, rather than to biological differences. More generally, no study on the repeatability of
416 provisioning behaviour has accounted for (i.e. removed) this sampling error (see Khwaja et al., 2017, for
417 summary of studies), suggesting that these estimates are underestimations, and that the comparison among
418 these studies may be problematic. Note that we have assumed in these calculations that this stochastic
419 error is Poisson-distributed. As discussed above, this may not be the case. However, this survey still
420 demonstrates how wide ranging the current underestimation of effect size is. Furthermore, as refractory
421 periods may vary between studies, not accounting for the stochastic error makes comparison between
422 studies extremely problematic.

423 ANALYTICAL SOLUTIONS

424 Directly modelling stochastic error

425 In the majority of analyses (80% and 63% in our surveys of 2015/16 and 2022, respectively), behavioural
426 counts are analysed as a response variable. Stochastic error can be accounted for in such analyses by
427 using a generalised linear mixed model (GLMM) framework, specifying a Poisson error distribution.
428 Broadly GLMMs account for distribution specific variance (i.e. the Poisson distributed stochastic error)
429 seen on the observed level, and transform the data from the expected scale onto the 'latent' scale, using
430 a 'link' function (in this case typically the log function). The data is normally distributed on the latent
431 scale, fulfilling linear model assumptions (see Bolker et al., 2009; Nakagawa and Schielzeth, 2010;
432 de Villemereuil et al., 2016, for a helpful breakdown of these models). We present examples in the
433 Supplementary Material (S6).

434 By using GLMMs, the stochastic variation is specifically modelled, and ICC and R^2 estimates can be
435 made with or without the Poisson distributed stochastic error (i.e. by including it or not in the calculation
436 of total variance; Figure 4). Traditionally R^2 and ICC are calculated with stochastic error (Nakagawa
437 and Schielzeth, 2010, 2013). However, the overriding utility of this has recently been discussed, and it
438 has been suggested that under some circumstances it is more appropriate to measure such metrics on the
439 expected scale (i.e without stochastic error; see de Villemereuil et al., 2016, for a more general discussion
440 of this in the context of heritability, with reference to both Poisson and Binomial stochastic error). In the
441 case of variables such as provisioning rate (that are sampled in a fixed period of time), we would argue that
442 these metrics should be estimated without the inclusion of this stochastic error in the estimation of total
443 variance (Figure 4), as this is dependent on the sampling effort, and so the metrics are more biologically
444 meaningful on the expected rather than the observed scale. Note that in some cases the actual behavioural
445 count, rather than the rate, is more relevant to a researcher's question, for example correlating the total
446 number of feeds within a time period with the mass change of chicks within that same time period. In
447 this case, it would not be appropriate to account for the stochastic error (i.e. by following the methods
448 proposed in Nakagawa and Schielzeth 2010, 2013, see also Nakagawa et al. 2017a), as the number of

449 feeds rather than the underlying rate would be the variable of interest.

450 Typically, provisioning rate has been analysed assuming a Gaussian error distribution (i.e. in linear
451 (mixed) models (L(M)Ms); 77% and 70 % of analyses in our surveys of 2015/16 and 2022, respectively).
452 The implicit assumption in these analyses is that all the variance can be explained i.e. that there is no
453 stochastic error that will always remain unexplained (Nakagawa and Schielzeth, 2013; Nakagawa et al.,
454 2017a, Figure 4). There has been much recent debate about the relative merits of using Gaussian or
455 Poisson error distributions to model count data (O’Hara and Kotze, 2010; Ives, 2015; Warton et al.,
456 2016; Morrissey and Ruxton, 2020). Here we emphasise that this stochastic error needs to be accounted
457 for (which is possible using either method; see Supplementary Material S6). Our recommendation for
458 GLMMs is therefore specifically based on the explicit estimation of stochastic error in these models,
459 leading to an output that is more intuitive and easier to deal with. Using LMMs to estimate effect sizes,
460 and post hoc removing stochastic error, can also result in estimates that are not bounded by 0 and 1 (the
461 limits of ICC and R^2). Gaussian distributions are also unbounded, suggesting that observations below 0 are
462 possible. A common alternative is to assume a log-Gaussian distribution (i.e. through log-transformation
463 of counts). Although this is bounded above 0, 0 is not a value that can exist, whilst being a possible value
464 for observed data (O’Hara and Kotze, 2010).

465 The resulting Poisson GLMMs are highly likely to be overdispersed, as it is unlikely that a given set
466 of predictors will explain all underlying variation. To account for this (additive) overdispersion, models
467 should be run as mixed models, including an observation level random effect (Hinde, 1982). The estimate
468 of variance from this observation level random effect can be used as the estimate for overdispersion
469 (non-distribution specific) variance, analogous to the residual variance of a linear model. Note that some
470 software (e.g. MCMCglmm; Hadfield, 2010) explicitly does this by default. It is also possible to model
471 such overdispersion in other ways, e.g. by assuming a negative binomial error distribution. This can
472 be parameterised as a Poisson-gamma mixture distribution (rather than the Poisson-log normal mixture
473 typically used in Poisson GLMMs).

474 It is worth noting that the parameter estimates themselves (i.e. estimates of intercepts, slopes and
475 variance components) should not change between the different models (log-normal and Poisson). These
476 effect sizes only differ when standardised by total variance (which is how metrics are typically compared
477 between studies). In studies of repeatability, we therefore also advocate reporting CV_B - coefficient of
478 between individual variation (Holtmann et al. 2017; see also Dochtermann and Royauté 2019). CV_B
479 reflects a mean-standardised measure of the amount of between individual variation. It is independent of
480 the method of analysis and the degree of stochastic error and so can readily be compared between studies,
481 regardless of sampling effort. It is analogous to CV_A (the coefficient of additive genetic variation; Houle,
482 1992), which was proposed to address similar issues of comparing additive genetic variation between
483 studies. Both CV_A and CV_B were first proposed in the context of Gaussian traits, and their derivation for
484 other distributions is more complex. Recently, de Villemereuil et al. (2016) derived a formulation of CV_A
485 for non-Gaussian traits, which holds also for CV_B . From a Poisson GLMM, CV_B is equal to the standard
486 deviation of the between individual effects on the latent scale. A demonstration of the calculation of CV_B
487 is presented in Supplementary Material S6.

488 Our assumption here is that the stochastic error is Poisson distributed, which may not be the case if,
489 for example, substantial refractory periods exist. The reduction in stochastic error due to such refractory
490 periods is also predictable, and can be modelled with a Tweedie distribution (see Supplementary Material
491 S1) or Generalised or Conway-Maxwell Poisson distributions (Lynch et al., 2014). Alternatively, the
492 length of intervals between behaviours can be modelled, which we discuss further below. Note that, with
493 or without a refractory period, current methods still act to systematically underestimate effect sizes, as they
494 do not account for stochastic error. It should also be noted that when refractory periods are small, assuming
495 Poisson distributed error results in less bias than assuming no stochastic error (see Supplementary material
496 S1). The choice that researchers make in how to calculate effect sizes (i.e. whether to remove stochastic
497 error, and if so the magnitude of that error) should be clearly defined in studies, allowing assessment and
498 future recalculation of relevant effect sizes.

499 Figure 6 demonstrates the effect of using different methods in the estimation of R^2 on simulated data
500 (Supplementary Material S5). In Figure 6A, we can see that by not correcting for this stochastic error
501 (using linear models; black dots), R^2 would increase as the mean number of observed visits increases and
502 would be systematically underestimated, as is seen in real provisioning data (Figure 3E). Accounting
503 for this error using Poisson GLMMs, results in (predominantly) unbiased estimates of R^2 (except at low

504 expected CV and mean number of observed visits; Figure 6C). The precision of these models is also
505 affected by both the expected CV and the mean number of observed visits, with precision increasing as
506 they both increase (Figure 6D).

507 **Behavioural counts as predictors - Measurement error models**

508 Although stochastic error can easily be modelled when behavioural counts are the response variables,
509 commonly used statistical software do not allow for the inclusion of error in predictor variables (linear
510 models, for example, assume that there is no error in the predictors). Indeed, no analysis in our literature
511 survey (out of 84 in 2015/16 and 29 in 2022) accounted for this stochastic error when provisioning rate
512 was used as a predictor variable. In order to model this error, we can use a class of models, known as
513 measurement error or ‘error in variable’ models, that allow error in the predictor variables to be specified.
514 These models are, however, complex to implement at present, although can readily be created in software
515 such as Stan (Carpenter et al., 2017). Measurement error models act very similarly to GLMMs, by creating
516 a latent variable (i.e. expected provisioning rate) which is then used as a predictor in the main model.
517 Variation in observation time between different observations can therefore easily be accounted for, as can
518 variables that may differ between observations. For example, if a researcher wanted to analyse the effect
519 of provisioning rate on chick mass at fledging, but provisioning rate had been measured at different brood
520 ages or in different environmental conditions at different nests, this variation could easily be accounted for.
521 In Supplementary Material S6 we present a practical example of such models in Stan (see also Freckleton,
522 2011; Garamszegi, 2016; Ponzi et al., 2018; Dingemans et al., 2021, for examples of the consequences of
523 measurement error and how to deal with it in ecology and evolution).

524 In Figure 6B we demonstrate the effect of not correcting for this stochastic error when using be-
525 havioural rate as a predictor variable (black dots); as the mean number of observed visits increases, the
526 predictive power of the behaviour increases, but is systematically downwardly biased. Measurement error
527 models (red dots) account well for this Poisson error, but as with GLMMs, their precision is low when the
528 mean number of observed visits is low.

529 **Analyse number, rather than transforming to rate**

530 In many studies the observation periods vary. This is frequently corrected for by scaling the observed
531 number of events to create a rate (e.g. visits/hour). In the provisioning literature, many studies also correct
532 for brood size when calculating provisioning rate (e.g. visits/hour/chick). Indeed, only 29% and 40% of
533 studies in our literature survey (in 2015/16 and 2022 respectively) analysed their data as a count rather
534 than a rate. However, when modelling count data, the raw number of observed events should be used,
535 as transformations will create several problems. Firstly, because of the way the mean, stochastic error
536 and expected variance scale (discussed above), by transforming the data the correct amount of Poisson
537 variance cannot be directly estimated (Figure 4). Once the number of observed arrivals is transformed
538 to a different scale (e.g. the number of arrivals is standardised to arrivals/hour), the mean no longer
539 represents the stochastic variance. For example, if the 20 ± 7.5 arrivals observed in study A from Figure
540 4 (120 mins), is transformed to 10 ± 3.75 arrivals/hour, the expected CV (i.e. the amount of biological
541 variation) is calculated as 0.2 (instead of 0.3), as we are overestimating the amount of Poisson sampling
542 error. Conversely, the 5 ± 2.7 arrivals in observations from study C (30 mins) when transformed becomes
543 10 ± 5.4 arrivals/hour, with an expected CV of 0.438. Therefore, depending on the direction of the scaling
544 (i.e. making the mean smaller or larger), this would lead to the respective over- or underestimation of
545 Poisson variance, and so a corresponding over- or underestimation of effect sizes (see also Supplementary
546 Materials S5 and S6). Instead, variation in observation time can be accounted for by using a Poisson
547 ‘exposure’ model (Gelman and Hill, 2007), by including log observation period as an offset (a covariate
548 with the slope fixed to 1). Similarly, brood size should be corrected for by including it as a covariate in the
549 GLMM. Correcting provisioning rate for brood size (e.g. using visits/hour/chick) incurs further problems
550 associated with the use of ratios, as the relationship between brood size and provisioning rate is often not
551 linear (Raubenheimer, 1995; Nakagawa et al., 2017b), and spurious correlations are created when brood
552 size is included as a covariate in addition to being corrected for in the response variable (Kronmal, 1993).

553 **Modelling interval lengths**

554 Given some of the problems outlined above, it may seem more appealing to model the length of the
555 intervals between behaviours rather than the count of the behaviours. Modelling interval lengths instead
556 of counts is not a solution in itself, however. Clearly, there is an advantage to analysing interval lengths

557 when auxiliary interval level data exists (such as environmental variables measured at the level of the
558 interval), as this provides an opportunity to start to understand what processes contribute to the apparent
559 stochastic nature of these interval lengths. However, auxiliary data are typically collected on the level of
560 the observation and not the interval, which does not give more information to the analysis than analysing
561 counts. Without interval-level data, the mean interval length for an observation is essentially the variable
562 of interest, and this is clearly a simple re-parametrisation of the counts (mean interval length = observation
563 period / count). As discussed above, we know that the interval lengths within an observation will show a
564 high level of stochastic error (they are expected to be exponentially distributed under the Poisson process
565 model). They can be treated as repeated measurements, where the within observation variation in interval
566 lengths represents the stochastic error.

567 Given our extensive discussion above, we would encourage the use of a gamma distribution when
568 modelling behavioural interval lengths. This approach also allows a population level α to be estimated,
569 and so some of the assumptions of a Poisson process can be directly assessed. An interesting extension to
570 this, would be to model whether both the refractory period and the return rate vary between observations
571 and further whether they systematically differ between individuals or environments.

572 **DETERMINING SAMPLING EFFORT**

573 Researchers frequently consider sample size when planning studies. In this case, that would typically
574 refer to the total number of observations. However, as we have shown, the number of events that are
575 observed within each observation is also important. How then should we determine the best strategy for
576 collecting such data?

577 Readers may wonder why we do not simply recommend extremely short observation periods, given
578 that we can correct for the additional stochastic error that is induced by this method. It is important to
579 note, however, that when observing a low mean number of events, precision of model parameter estimates
580 is low, and bias (under certain conditions) is high (Figure 6). This is because when the mean number of
581 observed visits is low, the estimation of the residual, unexplained biological variation is poor. As the
582 variance is so dominated by stochastic error, small random fluctuations in the mean (induced by sampling
583 error) have disproportionately large effects on the estimation of residual variance and can even lead to the
584 mean being larger than the observed variance, implying that there is no variation in expected rates. We
585 can see this pattern in both simulated data and in data from the literature survey (Supplementary Material
586 S5 and Figure S3). Moreover, this is likely why we see an upward bias in effect size at low mean values
587 in the Poisson models, as the residual (i.e. expected) variance is underestimated.

588 We should therefore seek to collect data under the conditions which minimise such effects. Ideally
589 data would be collected in a fully automated way, meaning that all (or a large proportion of) events would
590 be recorded, and the stochastic error across observations would be negligible. However, this is overly
591 idealistic in most situations, as setting up such systems involves a large amount of time and money, and
592 requires a high proportion of the population to be tagged to be effective. Thus, we advise using existing
593 data (or a pilot study) to estimate a suitable observation period (see Supplementary Material S5 for how
594 to calculate expected CV). This is not a one size fits all situation - an appropriate observation period will
595 differ among study systems, according to the mean rate and variability of the behaviour. The emphasis,
596 therefore, should be on the optimal mean number of observed events rather than optimal observation
597 period, as it is the former that will directly determine the proportion of stochastic error. Our simulations
598 show that an average of 20 events per observation minimises bias and maximises precision (but note
599 that the results may vary according to parameters such as the simulated R^2). We recognise that longer
600 observations may limit the number of observations that can be made, although researchers can use tools
601 such as planned missing data designs (Noble and Nakagawa, 2021) to offset this cost.

602 Finally, it is worth noting that our calculations are made on the assumption that a Poisson process is
603 the most suitable model for the behavioural count data in question, which may not be the case (see above).
604 However, regardless of the exact form of the stochastic error, extending observation periods will act to
605 reduce this error, and make comparisons between studies more meaningful.

606 **A CAUTIONARY NOTE ON OTHER MEASURES OF BEHAVIOUR**

607 Whilst the literature survey presented here focused on studies that specifically measured and analysed
608 counts of provisioning, there were many studies in our original search that analysed variables that are

609 derived from visit rate, such as the amount of food brought to the nest or proportion of visits made by
610 each sex. These metrics will equally be affected by the problems caused by stochastic error, as they
611 depend on visit rate for their quantification, and therefore the stochastic error associated with visit rate is
612 propagated to these other variables. We, therefore, would urge a similar word of warning in the use of any
613 measures derived from short observations of behaviour. The kind of stochastic error we describe here
614 applies not only to counts, but to any quantification of behaviour sampled in a short period of time. These
615 problems can be resolved through careful thought about which distribution to use in the analysis, and the
616 assumptions that a distribution has. For example, the amount of food brought to the nest might be well
617 described by compound Poisson or Tweedie distributions ((see Thompson, 1984, for an example with
618 rainfall data).

619 CONCLUSIONS

- 620 1. Stochastic error arises when measuring behaviour by counting the frequency of events in a sample
621 period. The degree of this error depends on both the number of events observed and the variation
622 in rates between observations. By not taking this error into account, we limit both the variation
623 in these behaviours that we can explain and the utility of these variables as predictors of other
624 traits. Furthermore, as the degree of this error depends on characteristics of the study, comparisons
625 between studies are highly problematic.
- 626 2. Using the null model of a Poisson process to describe this stochastic error, we can demonstrate it
627 arises in a predictable manner, allowing researchers to account for it using established statistical
628 methods. Whether, and how far, real behaviour count data deviates from this model is not well
629 understood. Future work should seek to address this, as it will give a better biological understanding
630 of the respective behaviour.
- 631 3. Using the example of provisioning rate, we demonstrate the suitability of the Poisson process as
632 a null model. However, through a literature survey we show that by far the majority of studies
633 of provisioning rate do not account for this stochastic error and, due to the low mean number of
634 observations per study, the amount of stochastic error is high. Whilst recent work may be correct
635 in asserting that provisioning rate is not an accurate descriptor of parental investment (Williams,
636 2012; Williams and a. Fowler, 2015), the methods that are currently employed to assess this are
637 insufficient to draw this conclusion. Therefore, although we welcome further investigation of the
638 different ways parents invest in their offspring (e.g. size / quantity / quality of prey), we suggest
639 that we should not yet rule out the possibility that provisioning rate itself, when properly measured,
640 is an adequate description of postnatal parental investment. The use of longer observations and
641 correct statistical analysis will aid us in these endeavours.
- 642 4. Finally, given the inevitability and predictable nature of this stochastic error, we should endeavour
643 to quantify and account for it when analysing behavioural count data, as well as taking steps to
644 minimise it where possible. Behavioural ecology is a discipline already fraught with relatively
645 small effect sizes, and low power to detect them; we do ourselves a disservice by adding more error
646 into the equation.

647 ACKNOWLEDGEMENTS

648 We thank the LHB discussion group at the University of Sheffield for valuable input, and Daniel Noble,
649 Fonti Kar, Alison Pick, David Westneat, Jarrod Hadfield, Paul Johnson and an anonymous reviewer for
650 their comments on the manuscript.

651 DATA AVAILABILITY

652 All data and code for analyses and simulations presented here can be found at <https://doi.org/10.5281/zenodo.7439115>

653 REFERENCES

654 Beekman, M. and Jordan, L. A. (2017). Does the field of animal personality provide any new insights for
655 behavioral ecology? *Behavioral Ecology*, 28(3):617–623.

- 656 Bell, A. M., Hankison, S. J., and Laskowski, K. L. (2009). The repeatability of behaviour: A meta-analysis.
657 *Animal Behaviour*, 77(4):771–783.
- 658 Blackwell, P. G., Niu, M., Lambert, M. S., and Lapoint, S. D. (2016). Exact Bayesian inference for animal
659 movement in continuous time. *Methods in Ecology and Evolution*, 7(2):184–195.
- 660 Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H. H., and White, J.
661 S. S. (2009). Generalized linear mixed models: A practical guide for ecology and evolution. *Trends in*
662 *Ecology and Evolution*, 24(3):127–135.
- 663 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo,
664 J., Li, P., and Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of Statistical*
665 *Software*, 76:1–32.
- 666 Daley, D. and Vere-Jones, D. (2003). *An Introduction to the Theory of Point Processes: Volume I:*
667 *Elementary Theory and Methods*. Springer-Verlag, New York.
- 668 de Villemereuil, P., Schielzeth, H., Nakagawa, S., and Morrissey, M. (2016). General methods for
669 evolutionary quantitative genetic inference from generalized mixed models. *Genetics*, 204(3):1281–
670 1294.
- 671 Dingemans, N. J., Araya-Ajoy, Y. G., and Westneat, D. F. (2021). Most published selection gradients are
672 underestimated: Why this is and how to fix it. *Evolution*, 75:806–818.
- 673 Dochtermann, N. A. and Royauté, R. (2019). The mean matters: Going beyond repeatability to interpret
674 behavioural variation. *Animal Behaviour*, 153:147–150.
- 675 Freckleton, R. P. (2011). Dealing with collinearity in behavioural and ecological data: Model averaging
676 and the problems of measurement error. *Behavioral Ecology and Sociobiology*, 65(1):91–101.
- 677 Garamszegi, L. Z. (2016). A simple statistical guide for the analysis of behaviour when data are constrained
678 due to practical or ethical reasons. *Animal Behaviour*, 120:223–234.
- 679 Gelman, A. and Hill, J. (2007). *Data Analysis Using Regression and Multilevel Hierarchical Models*.
680 Cambridge University Press, Cambridge.
- 681 Hadfield, J. D. (2010). MCMC methods for multi-response generalized linear mixed models: The
682 {MCMCglmm} {R} package. *Journal of Statistical Software*, 33(2):1–22.
- 683 Harrison, F., Barta, Z., Cuthill, I., and Székely, T. (2009). How is sexual conflict over parental care
684 resolved? A meta-analysis. *Journal of Evolutionary Biology*, 22(9):1800–1812.
- 685 Heuer, A., Mueller, C., and Rubner, O. (2010). Soccer: Is scoring goals a predictable Poissonian process?
686 *Europhysics Letters*, 89:38007.
- 687 Hinde, J. (1982). Compound Poisson Regression Models. In Gilchrist, R., editor, *GLIM 82: Proceedings*
688 *of the International Conference on Generalised Linear Models*, pages 109–121. Springer New York.
- 689 Holtmann, B., Lagisz, M., and Nakagawa, S. (2017). Metabolic rates, and not hormone levels, are a
690 likely mediator of between-individual differences in behaviour: A meta-analysis. *Functional Ecology*,
691 31:685–696.
- 692 Houle, D. (1992). Comparing evolvability and variability of quantitative traits. *Genetics*, 130:195–204.
- 693 Ihle, M., Pick, J. L., Winney, I. S., Nakagawa, S., and Burke, T. (2019). Measuring Up to Reality:
694 Null Models and Analysis Simulations to Study Parental Coordination Over Provisioning Offspring.
695 *Frontiers in Ecology and Evolution*, 7:142.
- 696 Ives, A. R. (2015). For testing the significance of regression coefficients, go ahead and log-transform
697 count data. *Methods in Ecology and Evolution*, 6(7):828–835.
- 698 Johnstone, R. A., Manica, A., Fayet, A. L., Stoddard, M. C., Rodriguez-Gironés, M. A., and Hinde,
699 C. A. (2014). Reciprocity and conditional cooperation between great tit parents. *Behavioral Ecology*,
700 25(1):216–222.
- 701 Khwaja, N., Preston, S. A. J., Hatchwell, B. J., Briskie, J. V., Winney, I. S., and Savage, J. L. (2017).
702 Flexibility but no turn-taking in provisioning riflemen (*Acanthisitta chloris*). *Animal Behaviour*,
703 125:25–31.
- 704 Kiureghian, A. D. and Ditlevsen, O. (2009). Aleatory or epistemic? Does it matter? *Structural Safety*,
705 31(2):105–112.
- 706 Koenig, W. D. and Walters, E. L. (2016). Provisioning patterns in the cooperatively breeding acorn
707 woodpecker: Does feeding behaviour serve as a signal? *Animal Behaviour*, 119:125–134.
- 708 Kronmal, R. A. (1993). Spurious Correlation and the Fallacy of the Ratio Standard Revisited. *Journal of*
709 *the Royal Statistical Society*, 156(3):379–392.
- 710 Lendvai, A. Z., Akçay, C., Ouyang, J. Q., Dakin, R., Domalik, A. D., St. John, P. S., Stanback, M., Moore,

- 711 I. T., and Bonier, F. (2015). Analysis of the optimal duration of behavioral observations based on an
712 automated continuous monitoring system in tree swallows (*Tachycineta bicolor*): Is one hour good
713 enough? *PLoS ONE*, 10(11):1–11.
- 714 Levin, R. J. (2009). Revisiting Post-Ejaculation Refractory Time—What We Know and What We Do Not
715 Know in Males and in Females. *The Journal of Sexual Medicine*, 6(9):2376–2389.
- 716 Lynch, H. J., Thorson, J. T., and Shelton, A. O. (2014). Dealing with under- and over-dispersed count
717 data in life history, spatial, and community ecology. *Ecology*, 95(11):3173–3180.
- 718 Morrissey, M. B. and Ruxton, G. D. (2020). Revisiting advice on the analysis of count data. *Methods in
719 Ecology and Evolution*, 11(9):1133–1140.
- 720 Morvai, B., Nanuru, S., Mul, D., Kusche, N., Milne, G., Székely, T., Komdeur, J., Miklósi, A., and Pogány,
721 A. (2016). Diurnal and Reproductive Stage-Dependent Variation of Parental Behaviour in Captive
722 Zebra Finches. *PLoS ONE*, 11(12):e0167368.
- 723 Murphy, M. T., Chutter, C. M., and Redmond, L. J. (2015). Quantification of avian parental behavior:
724 What are the minimum necessary sample times? *Journal of Field Ornithology*, 86(1):41–50.
- 725 Nakagawa, S. and Cuthill, I. C. (2007). Effect size, confidence interval and statistical significance: A
726 practical guide for biologists. *Biological Reviews of the Cambridge Philosophical Society*, 82:591–605.
- 727 Nakagawa, S., Johnson, P. D., and Schielzeth, H. (2017a). The coefficient of determination R² and
728 intra-class correlation coefficient from generalized linear mixed-effects models revisited and expanded.
729 *Journal of the Royal Society Interface*, 14:20170213.
- 730 Nakagawa, S., Kar, F., O’Dea, R. E., Pick, J. L., and Lagisz, M. (2017b). Divide and conquer? Size
731 adjustment with allometry and intermediate outcomes. *BMC Biology*, 15(1):107.
- 732 Nakagawa, S. and Schielzeth, H. (2010). Repeatability for Gaussian and non-Gaussian data: A practical
733 guide for biologists. *Biological Reviews of the Cambridge Philosophical Society*, 85(4):935–56.
- 734 Nakagawa, S. and Schielzeth, H. (2013). A general and simple method for obtaining R² from generalized
735 linear mixed-effects models. *Methods in Ecology and Evolution*, 4:133–142.
- 736 Noble, D. W. A. and Nakagawa, S. (2021). Planned missing data designs and methods: Options for
737 strengthening inference, increasing research efficiency and improving animal welfare in ecological and
738 evolutionary research. *Evolutionary Applications*, 14(8):1958–1968.
- 739 O’Hara, R. B. and Kotze, D. J. (2010). Do not log-transform count data. *Methods in Ecology and
740 Evolution*, 1(2):118–122.
- 741 Pick, J. L., Nakagawa, S., and Noble, D. W. (2019). Reproducible, flexible and high throughput data
742 extraction from primary literature: The metaDigitise R package. *Methods in Ecology and Evolution*,
743 10:426–431.
- 744 Poisot, T. (2011). The digitize package: extracting numerical data from scatterplots. *The R Journal*,
745 3(1):25–26.
- 746 Ponzi, E., Keller, L. F., Bonnet, T., and Muff, S. (2018). Heritability, selection, and the response to selection
747 in the presence of phenotypic measurement error: Effects, cures, and the role of repeated measurements.
748 *Evolution*, 72(10):1992–2004. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/evo.13573>.
- 749 Raubenheimer, D. (1995). Problems with Ratio Analysis in Nutritional Studies. *Functional Ecology*,
750 9(1):21–29.
- 751 Savage, J. L., Browning, L. E., Manica, A., Russell, A. F., and Johnstone, R. A. (2017). Turn-taking in
752 cooperative offspring care: By-product of individual provisioning behavior or active response rule?
753 *Behavioral Ecology and Sociobiology*, 71:162.
- 754 Schlicht, E., Santema, P., Schlicht, R., and Kempnaers, B. (2016). Evidence for conditional cooperation
755 in biparental care systems? A comment on Johnstone et al. *Behavioral Ecology*, 27:e2–e5.
- 756 Schwagmeyer, P. L. and Mock, D. W. (2008). Parental provisioning and offspring fitness: Size matters.
757 *Animal Behaviour*, 75(1):291–298.
- 758 Spence, M. A., Muiruri, E. W., Maxwell, D. L., Davis, S., and Sheahan, D. (2021). The application of
759 continuous-time markov chain models in the analysis of choice flume experiments. *Journal of the
760 Royal Statistical Society: Series C (Applied Statistics)*, 70(4):1103–1123.
- 761 Sánchez-Tójar, A., Schroeder, J., and Farine, D. R. (2018). A practical guide for inferring reliable
762 dominance hierarchies and estimating their uncertainty. *Journal of Animal Ecology*, 87(3):594–608.
- 763 Thompson, C. S. (1984). Homogeneity analysis of rainfall series: An application of the use of a realistic
764 rainfall model. *Journal of Climatology*, 4:609–619.
- 765 Wan, X., Wang, W., Liu, J., and Tong, T. (2014). Estimating the sample mean and standard deviation

766 from the sample size, median, range and/or interquartile range. *BMC medical research methodology*,
767 14(1):135.

768 Warton, D. I., Lyons, M., Stoklosa, J., Ives, A. R., and Schielzeth, H. (2016). Three points to consider
769 when choosing a LM or GLM test for count data. *Methods in Ecology and Evolution*, 7(8):882–890.

770 Westneat, D. F., Mutzel, A., Bonner, S., and Wright, J. (2017). Experimental manipulation of brood
771 size affects several levels of phenotypic variance in offspring and parent pied flycatchers. *Behavioral
772 Ecology and Sociobiology*, 71(6).

773 Westneat, D. F., Schofield, M., and Wright, J. (2013). Parental behavior exhibits among-individual
774 variance, plasticity, and heterogeneous residual variance. *Behavioral Ecology*, 24:598–604.

775 Williams, T. D. (2012). *Physiological Adaptations for Breeding in Birds*. Princeton University Press,
776 Princeton.

777 Williams, T. D. and a. Fowler, M. (2015). Individual variation in workload during parental care: Can
778 we detect a physiological signature of quality or cost of reproduction? *Journal of Ornithology*,
779 156:441–451.

780 Wolak, M. E., Fairbairn, D. J., and Paulsen, Y. R. (2012). Guidelines for estimating repeatability. *Methods
781 in Ecology and Evolution*, 3(1):129–137.

782 FIGURES AND TABLES

Table 1. Glossary of terms used in the manuscript

Point/event	Occurrence of a behaviour
Rate	Number of events per unit time
Observation	Observation of the unit of interest (e.g. an individual or a nest) for a defined period of time
Observation period (t)	The length of the observation (e.g. one hour)
Interval	The interval between two events (e.g. arrivals at a nest)
Interval length	The length of time between two events
Point process	Statistical description of events occurring through time
Poisson process	Simple point process with a single parameter, the rate (λ).
λ	‘True’ rate, an underlying/latent, unmeasurable variable. Equal to 1/expected interval length. When we use the number of arrivals in an observation period or the mean interval length, we are implicitly estimating this quantity.
Expected number of events	λt i.e. the number of events we would <i>expect</i> to see in a given time period and a given occurrence rate
Observed number of events (y)	Number of events actually <i>observed</i> in an observation period
Stochastic error	Error induced in our estimate of rate through sampling design
Refractory period	A period in which a behaviour is unlikely to reoccur

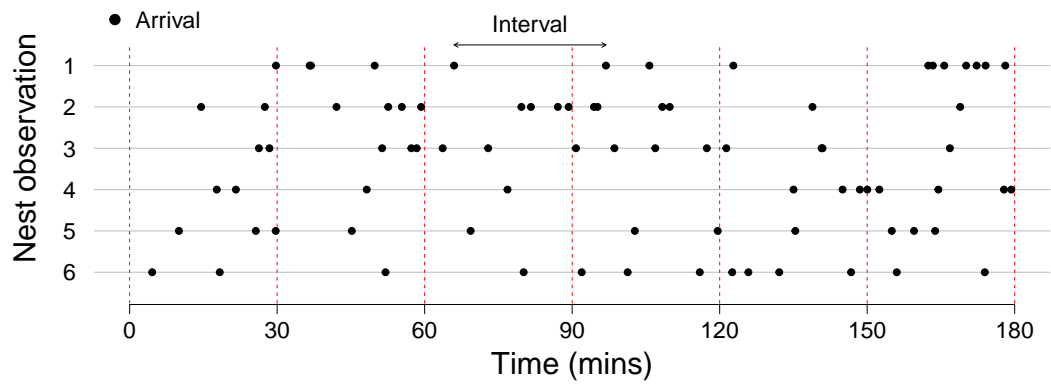


Figure 1. Behaviour can be described as points (or events) occurring on a straight line (through time), in other words as a point process. 6 nest observations are shown (grey lines), that were all simulated using a Poisson process with the *same* arrival rate (4.5 arrivals/hour) to demonstrate the variation that may arise through observations of different nests with the same arrival rate. The red dotted lines demonstrate the effect of shortening observation periods on the variation between observations.

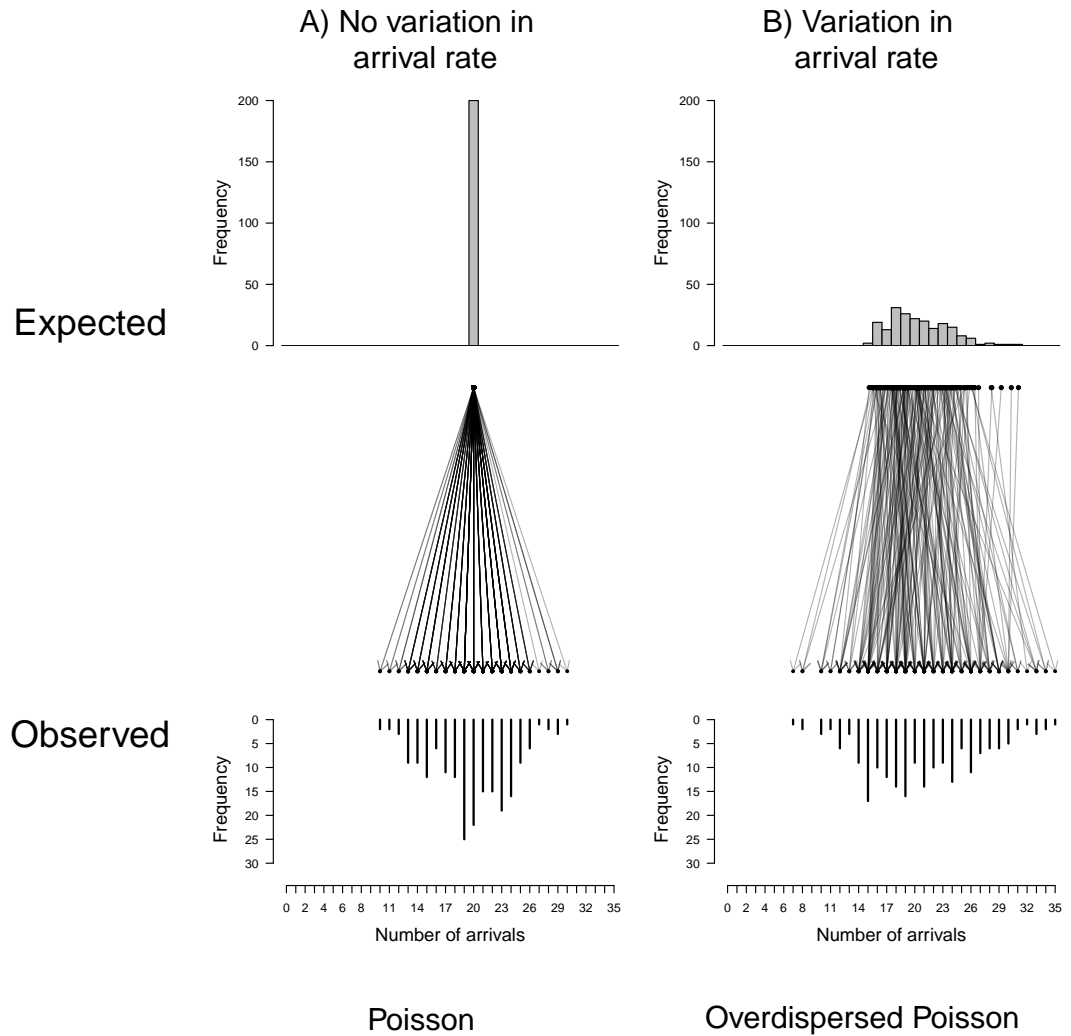


Figure 2. Visualisation of Poisson distributed stochastic error. A) If all observations are the same length and have the same expected rate (e.g. $t = 2$, $\lambda = 10$ and $\sigma_{\lambda t}^2 = 0$, so $\lambda t = 20$), the number of visits across all observations would be Poisson distributed ($\sigma_{stoc}^2 = 20$). B) When there is variation in the expected rate (for example, due to consistent differences between individuals; $\sigma_{\lambda t}^2 > 0$), every different rate is observed with stochastic error, leading to an over-dispersed Poisson distribution on the observed scale.

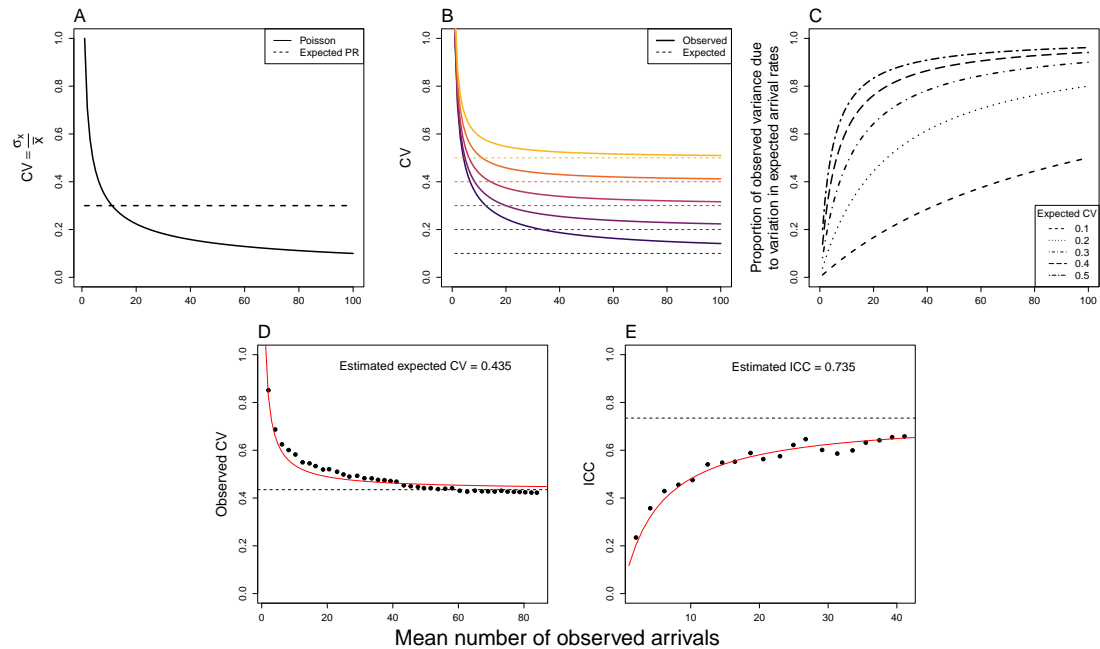


Figure 3. In data with Poisson distributed stochastic error, we expect to see certain patterns. Because the coefficient of variation (CV) of Poisson distributed stochastic error (solid line) and expected arrival rate (dotted line), change differently as the mean number of observed arrivals increases (A), we see that the observed CV (solid line) decreases as the mean number of observed arrivals increases, according to the variation in expected rates (B). In D) we can see this pattern in real provisioning data. As a result, the proportion of total observed variation due to ‘biological’ variation in expected rates, should increase with the mean number of observed arrivals (C). This proportion can be interpreted as either the maximum repeatability or R^2 when behavioural counts are analysed as a response variable or the maximum amount of variation these counts can explain in another variable, when Poisson distributed stochastic error is not accounted for. In E), we see this pattern arising in real provisioning data; repeatability (ICC) increases with increasing observation period. D) and E) use data presented in Lendvai et al. (2015); the red line shows the predictions from a non-linear model estimating CV and ICC, respectively, assuming an underlying Poisson process, and dotted line the estimated CV and ICC from these models (see Supplementary Material S2).

Example Population

Arrival Rate = 10 visits/hour, Expected CV = 0.3, Repeatability (ICC) = 0.5

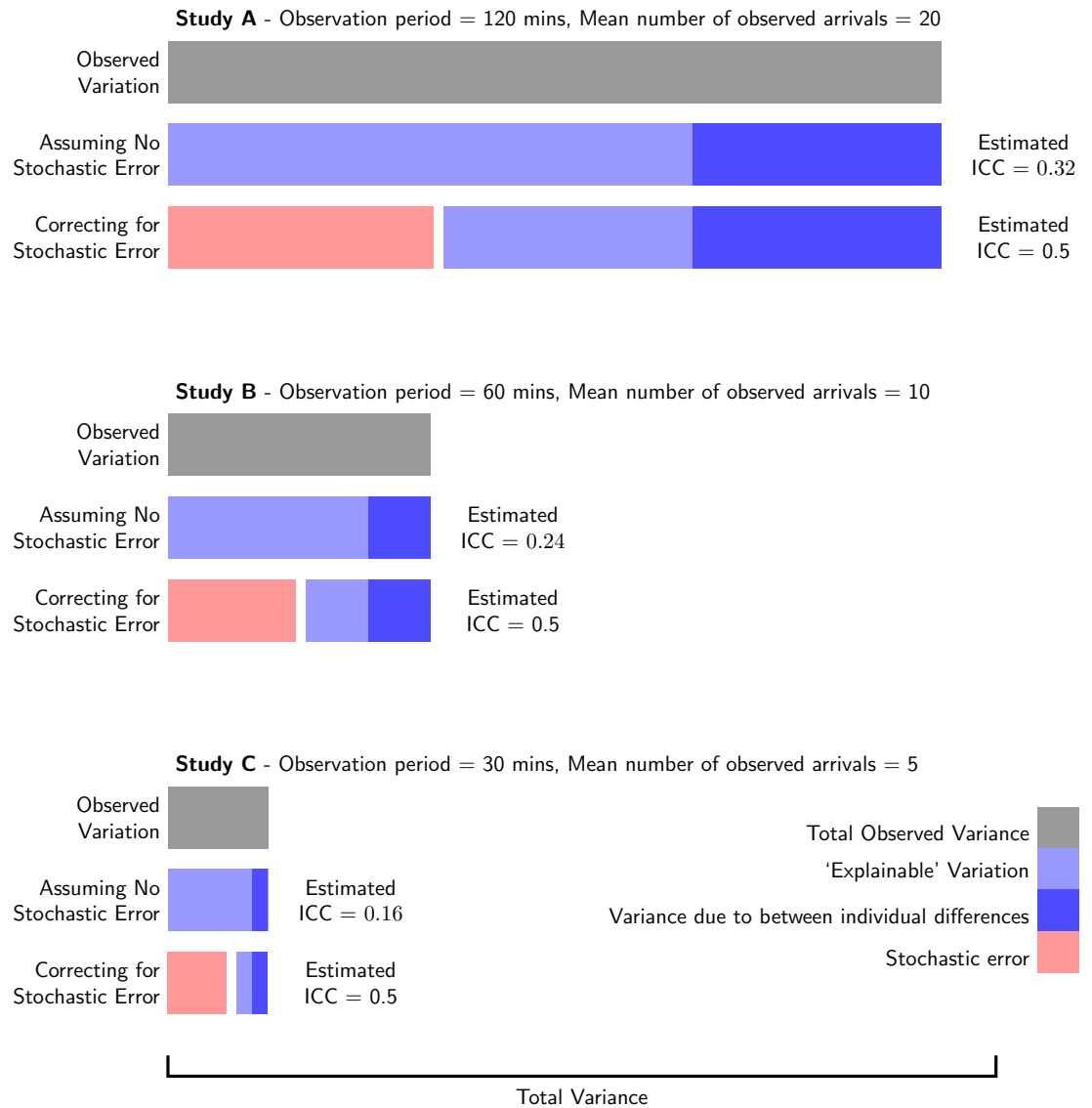


Figure 4. Effect of not accounting for Poisson distributed stochastic error in analyses of behavioural count data. Three studies of different observation periods in the same population will estimate different effect sizes, when not accounting correctly for the presence of stochastic error.

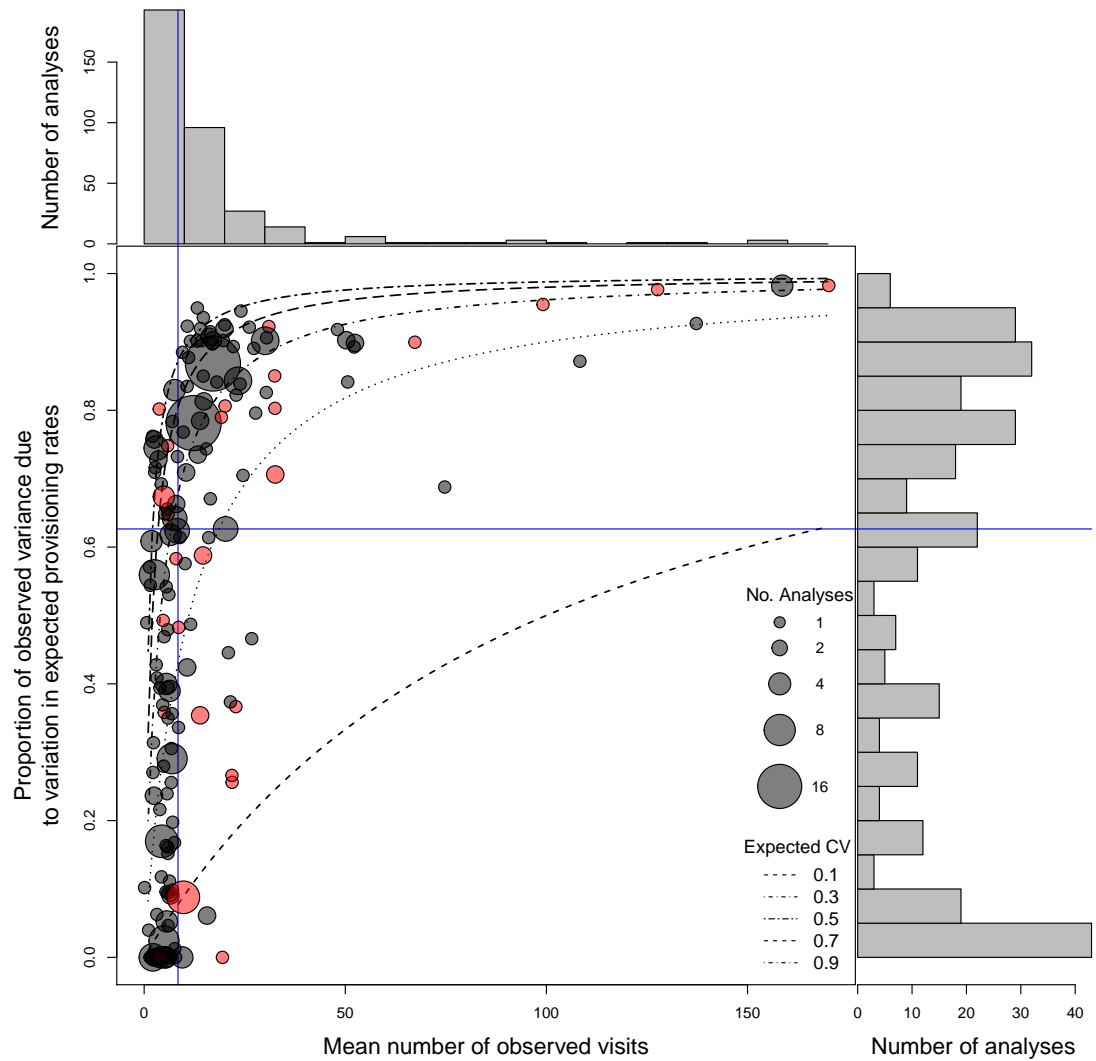


Figure 5. Distributions of the mean number of observed visits and proportion of observed variation due to expected variation in provisioning rates, from provisioning data used in analyses presented in papers published in 2015/16, and the relationship between them. Points show individual datasets, the size of the points indicates the number of models that were run using this dataset. Grey points indicate data from direct observation or video recordings, and red from automated data collection. Blue lines show median values. Estimates for which the mean was greater than the variance, the proportion of observed variation is displayed as 0.

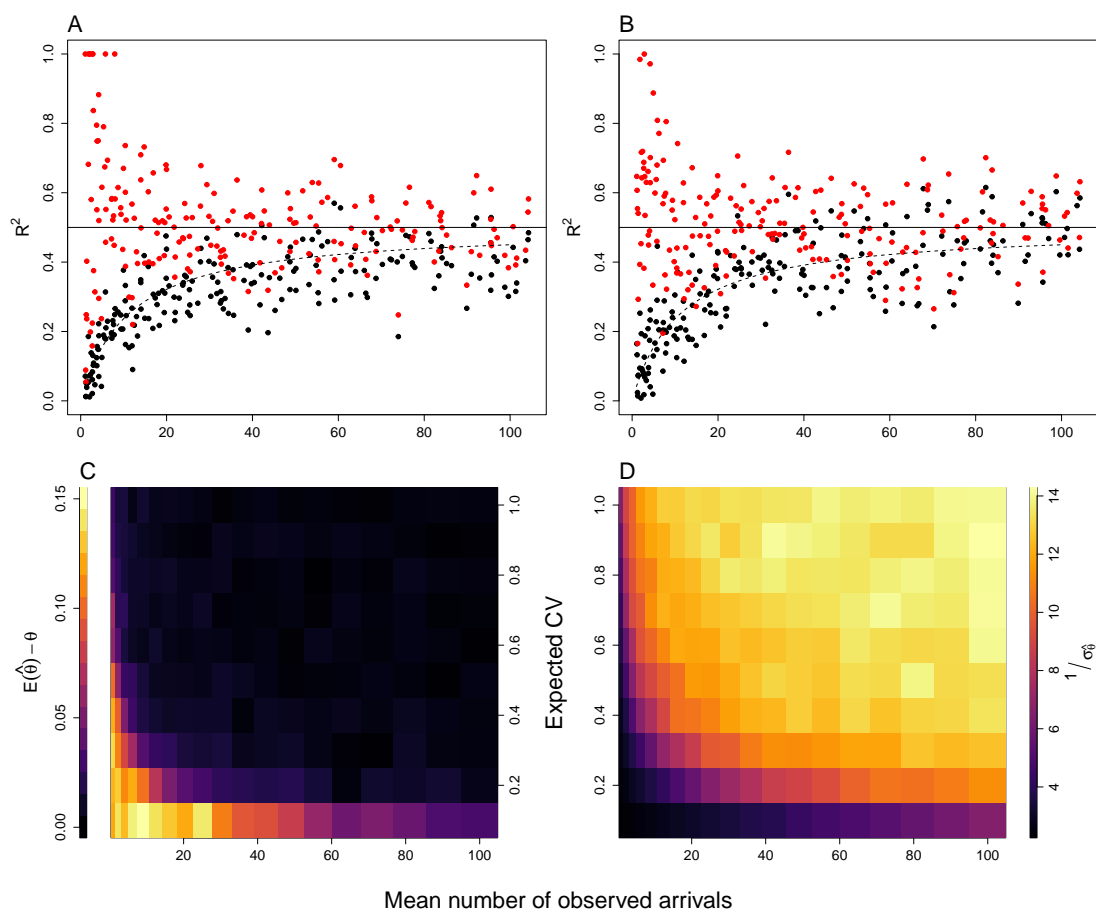


Figure 6. Results of simulations showing the effect of analysing behavioural count data with (red points) and without (black points) accounting for Poisson error, as both response (A) and predictor (B), over varying mean number of observed visits. A) and B) were simulated with an expected CV of 0.3; solid lines show simulated R^2 , and dotted lines show predicted R^2 when not taking in account Poisson error. C) and D) show bias and precision, respectively, in R^2 calculated from Poisson GLMM with provisioning rate as a response variable, from simulations across varying means and expected CVs; is the simulated value, and is the estimated value. Blue colours show low bias and low precision, respectively, and red colours high bias and high precision. See Supplementary Material S5 for further details of the simulations.