# Dispersion tests in generalised linear mixed-effects models - a methods comparison and practical guide for ecologists

Melina de Souza Leite[1*], Daniel Rettelbach[1,2] & Florian Hartig[1]

1. Theoretical Ecology, University of Regensburg, Germany
2. coTrial Associates, Department of Surgery, University Hospital Regensburg, Germany (current address)

*corresponding author: melina.souza-leite@ur.de

**Headline:** Dispersion tests for GLMMs

## Author contributions

MSL, FH, and DR conceived the ideas and designed the methodology. MSL wrote the simulation code and created the final version of the graphs and tables. MSL and FH led the writing of the manuscript. All authors contributed critically to the drafts and gave final approval for publication.

## Data availability

All data were simulated. The code for the simulations, analyses, and figures is available on Zenodo (https://doi.org/10.5281/zenodo.17611061).

## Acknowledgements

## Conflict of interest

21    The authors, MSL and FH, are developers of the R package DHARMa, which

22    implements the dispersion tests used in this study.

## Abstract

24    1.  Underdispersion and overdispersion are common issues when analysing

25        ecological data with generalised linear (mixed) models (GLMs/GLMMs).

26        Overdispersion, the phenomenon where observations spread wider than expected

27        by the fitted model, usually leads to anti-conservative p-values and, thus, to

28        inflated type I error. In contrast, underdispersion, a narrower spread of the data

29        than expected, causes overly conservative p-values and, therefore, reduced

30        power. A range of tests has been proposed to detect such dispersion problems,

31        but there are few comparative studies of their performance across models and

32        analysis settings.

33    2.  The goal of this study is to identify a general dispersion test for GLMs/GLMMs

34        that is applicable across all standard distributions and random-effects structures.

35        Following an initial assessment of available tests, we selected two classes of

36        dispersion tests as candidates: (1) parametric and nonparametric tests based on

37        Pearson residuals and (2) simulation-based tests that compare the expected and

38        observed residual variance.

39    3.  Comparing their performance by type I error, power, and dispersion estimate,

40        across a range of GLMs and GLMMs, we found that the nonparametric Pearson

41        residuals test performed best across all metrics, especially for data with low

42        incidence or count rates and/or small samples; however, at the cost of high

43        computational expense. The parametric Pearson residuals test, recommended in

44        many books and guidelines, was fast and effective for GLMs, but biased towards

45      underdispersion in GLMMs due to the naïve computation of the random-effect

46      degrees of freedom. The simulation-based residual variance test was slightly less

47      powerful, but showed overall good calibration and was much faster to compute.

48      The latter offers a compromise between the strengths and weaknesses of the two

49      Pearson-based tests.

50      4.  We conclude that for GLMs, the parametric Pearson residuals test offers the best

51      balance of speed and accuracy. For GLMMs, we recommend either the

52      computationally demanding nonparametric Pearson residuals test or the faster,

53      although somewhat less powerful, simulation-based residual variance test. We

54      also provide additional recommendations for ecological data analysis to address

55      dispersion issues using the most commonly used R packages, avoiding pitfalls

56      and improving model fit and the interpretation of ecological datasets.

57      **Keywords:** overdispersion/underdispersion, multilevel/hierarchical models, hypothesis

58      test, Pearson residuals, type I error, power, dispersion parameter

59      **Code for Peer review:**

60      https://anonymous.4open.science/r/dispersion_test_GLMM/README.md

## Introduction

61

62 Generalised linear models (GLMs) and generalised linear mixed models (GLMMs) are

63 the most commonly used tools for the statistical analysis of ecological data (Bolker et

64 al., 2009; Lai et al., 2019; Touchon & McCoy, 2016). By incorporating mixed and

65 random effect structures with a wide array of distributional assumptions (e.g., binomial,

66 Poisson), GLMMs allow researchers to model nonnormal response variables (e.g.,

67 counts, proportions, or presence-absence) while properly accounting for variation

68 clustered in sampling units, sites, or study years (Bolker et al., 2009; McMahon & Diez,

69 2007). However, as for all parametric statistics, these models rely on the fact that

70 residuals scatter around the regression mean with the specified distribution, and their

71 inferential results can be seriously biased if these distributional assumptions are

72 violated.

73        A particularly common and dreaded violation of distributional assumptions in

74 GLMs/GLMMs is overdispersion. Overdispersion refers to greater variation in the

75 observed data (and particularly the model residuals) than the fitted model assumes

76 (Campbell, 2021; McCullagh & Nelder, 1989). Strong overdispersion usually appears in

77 distributions that assume a fixed mean-variance relationship, such as the Poisson model

78 for count data (Harrison, 2014; Hilbe, 2014) or the binomial model for discrete

79 proportions (Dunn & Smyth, 2018; Harrison, 2015). For example, a Poisson process

80 assumes that we count randomly distributed points in space. However, when individuals

81 are subject to spatial/temporal clustering due to different ecological mechanisms (e.g.,

82 patchy resource distribution, social behaviour, dispersal limitation) and/or imperfect

83 detection (Rhodes, 2015), we typically find higher dispersion than expected from a

84 Poisson distribution (Box 1). Alternatively, overdispersion may also arise from

85    modelling misfit, for example, by failing to include important predictors and

86    interactions or by specifying the incorrect link function (Hilbe, 2011).

87        Overdispersion is a major concern in practical data analyses because it can have

88    substantial anti-conservative effects on p-values, confidence intervals, and all other

89    goodness-of-fit and precision metrics (Fig. 1, see also Rhodes, 2015). Anti-conservatism

90    means that p-values and confidence intervals are too small, leading to inflated false-

91    positive results (type I errors). In practice, we have encountered analyses where an

92    overdispersed model had very small and significant p-values (<0.001) that became

93    nonsignificant after switching to a GLM with more appropriate dispersion (see example

94    in Fig. 1).

95        The counterpart to overdispersion is underdispersion, where the variation in the

96    observed data (and, thus, model residuals) is lower than assumed by the fitted model.

97    Reasons for underdispersion can again be that the data-generating process (e.g., a

98    uniform distribution of individuals in space, Box 1) differs from what is assumed by the

99    model (Lynch et al., 2014). However, in practice, it is often the result of model

100   overfitting, i.e., having a too complex model that overfits the data. Underdispersion is

101   somewhat less discussed in the ecological literature, both because it is less frequent, but

102   also because it leads to over-conservative model metrics (Fig. 1). This may seem less

103   problematic as it does not lead to reporting "wrong" effects, but underdispersion

104   reduces overall power and thus increases type II error. Therefore, accurate statistical

105   inference demands that we identify and adequately address both underdispersion and

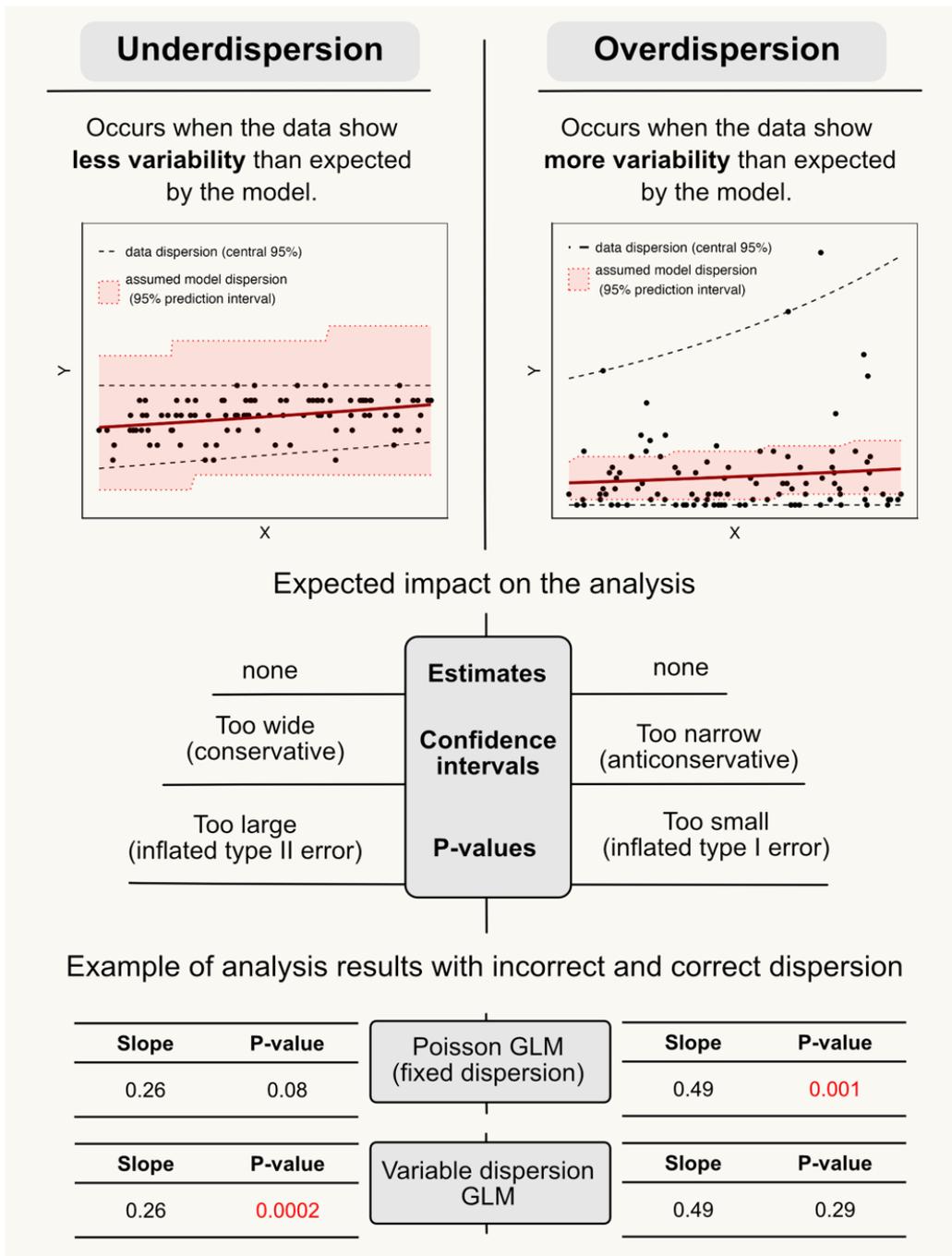106   overdispersion to minimise the risk of wrong inference.

107       Given the central importance of dispersion to all statistical indicators,

108   statisticians have pondered how to detect and address dispersion problems since the

109   early days of modern statistics in the 19[th] century (see Quine & Seneta (1987) and

110    Xekalaki (2014) for a historical perspective). Since then, a large variety of approaches

111    have been proposed and discussed to deal with the "dispersion problem", ranging from

112    (1) comparing models with or without free dispersion parameters through likelihood

113    ratio test, such as Poisson and negative binomial (e.g. Yang et al., 2007), (2) designing

114    specific hypothesis tests for the "extra" variation (e.g. Fisher, 1950), such as score tests

115    (Dean, 1992; Dean & Lawless, 1989; Lawless, 1987), (3) using goodness-of-fit tests,

116    such as tests on Pearson or Deviance residuals (Dunn & Smyth, 2018; McCullagh,

117    1985) (although the distinction between categories (2) and (3) can be blurry, see

118    (Collings & Margolin, 1985; Dean, 1992; Dean & Lawless, 1989) or (4) using

119    simulation-based non-parametric tests to compare observed and predicted variance of

120    the residuals (Hartig, 2024).

121        Somewhat confusing for the ecological data analyst, however, many of these

122    approaches have been designed and tested only in very specific scenarios (e.g. only for

123    a Poisson GLM), and there is a surprising lack of systematic evaluation of these tests

124    and strategies across a range of more complex GLMMs. Moreover, a quick review of

125    current methods in the R environment (R Core Team, 2024) revealed that existing

126    dispersion tests are scattered across different packages (Table 1), and most of them

127    work only for a restricted set of models. In the ecological literature, although awareness

128    of dispersion problems has increased over the last 20 years (Box 2), there is still a clear

129    lack of guidance on how dispersion problems are assessed or tested (Box 2). All this

130    makes it challenging to decide which test to use in applied ecological data analysis.

131        The goals of this study are: (1) to review and order the diversity of dispersion

132    tests for GLMs and GLMMs, and (2) to identify tests that can reliably work across a

133    range of models with diverse distributions and complex hierarchical structures, common

134    situations for ecological data analyses. Based on our literature review (next section), we

135     identified two groups of tests that appeared to be generally applicable: parametric and

136     non-parametric tests on Pearson residuals, as well as a new simulation-based non-

137     parametric test that directly compares observed and predicted variance of the raw

138     residuals. We then used simulated data to compare the performance of these tests in

139     terms of type I error, power, and the interpretability of the dispersion statistics. Based on

140     our results, we provide recommendations for the most suitable tests for detecting over-

141     or underdispersion, depending on model complexity and software availability, i.e.,

142     currently available R packages and functions.

**Figure 1.** Definition, statistical consequences, and a practical analysis example of under-/overdispersion in GLM/GLMMs. The top row shows examples of a data analysis using a Poisson GLM with simulated under- and overdispersed count data. Data points in black are contrasted to the Poisson model's 95% prediction interval (in red). Black dashed lines illustrate the data dispersion (central 95% quantiles). In the example, we present slope estimates and p-values for the GLM Poisson models fitted to the under- and overdispersed data above, and the results from more appropriate models with correct dispersion: a Conway-Maxwell-Poisson GLM for underdispersed data and a negative binomial GLM for overdispersed data.

**BOX 1: Ecological causes of under- and overdispersion**

Under/overdispersion in ecological data may not be only a statistical problem that we need to control for, but may also reflects key ecological processes of interest (Nakagawa et al., 2026; Rhodes, 2015). For example, ecological field data often consists of individual counts in space or time. If individuals are distributed completely at random, the sampling variability will follow a Poisson distribution. However, a range of ecological, observational and modelling processes can lead to deviations from this distribution, resulting in over- or underdispersion. For example, spatial aggregation (clustering) of individuals due to patchy resources distribution, social behaviour, or dispersal limitation increases sampling variability and thus creates overdispersion compared to the Poisson (Fig. B1, see also Lindén & Mäntyniemi, 2011). On the contrary, a uniform spatiotemporal distribution of individuals, for example due to territoriality, may create underdispersion (Lynch et al., 2014). Note that in these examples, but also in general, over- and underdispersion are always defined with respect to an expectation, in this case, the Poisson distribution.
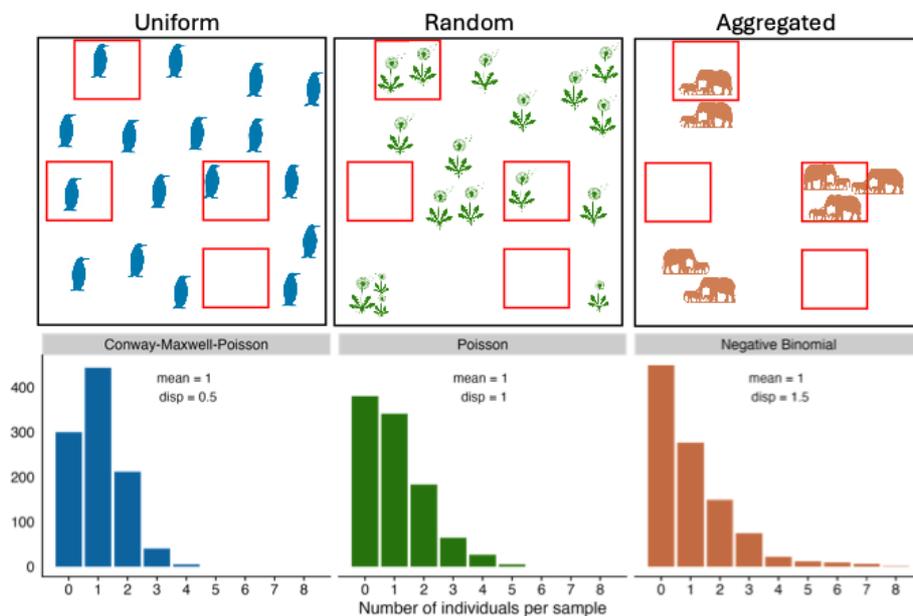


**Figure B1.** Examples of spatial distribution patterns of individuals under the same sampling design (quadrats in red) and the corresponding data-generating distributions for each pattern: Poisson (random pattern, green), negative binomial (aggregated pattern, light brown), and Conway-Maxwell-Poisson (uniform pattern, blue). Histograms are drawn from 1000 samples from the different distributions with the same mean (1) but varying dispersion (0.5, 1, 1.5). The figure silhouettes correspond to classical text-book examples (e.g., Campbell & Reece, 2005): most penguin species are territorials, tending to be uniformly spaced; dandelions have wind-dispersed seeds and tend to have a random distribution; elephants live in groups, and therefore exhibit an aggregated distribution.

More ecological causes for underdispersion appear, for example, in individual reproductive metrics (e.g, such as clutch size or seeds per fruit) with discrete counts upper limited by behavioural or physiological constrains, such as ovule number, parental care or resources availability (Brooks et al., 2019; Lynch et al., 2014; Puig et al., 2024).

It is worth noting that it doesn't mean, however, that dispersion problems always hint to an interesting ecological process. As discussed in the main text, dispersion problems may also arise from observational errors, for example imperfect detection (Rhodes, 2015), or from model misfit. When detecting dispersion problems, a careful consideration of their reasons is therefore paramount for an adequate ecological interpretation.

---

**BOX 2: Current practices for dispersion issues in ecological studies**

To understand the current practice for addressing dispersion problems in GLMs/GLMMs for count and discrete proportions data, we performed a text mining analysis of the ecological literature from the last 20 years (see S1 for details). Our results show that, over recent years, the percentage of all ecological studies using GLM/GLMMs for such data ahs remained around 8% (Fig. S1.1). Within these studies, we observed a steady increase in awareness about dispersion issues, with more than 28% of studies published in 2025 explicitly mentioning them (Fig. B2A).



**Figure B2. A)** Annual trends for the proportion of ecological studies using GLMs/GLMMs for count and/or discrete proportion data **that mention dispersion terms** in the text**. B)** The type of checks and tools used for dispersion problems found in the papers that mentioned dispersion terms in 100 randomly selected papers from 2025. For details of the text mining analysis see S1.

We further analysed a subset of 100 randomly selected ecological papers in 2025 that used GLMs/GLMMs and mentioned dispersion issues in more detail. 81 mentioned overdispersion, 4 underdispersion and 4 tested for both issues (see S1). Among them, only 40 papers explicitly reported testing for dispersion problems (Fig. B2B): 25 using the DHARMa package (but not mentioning which test), 5 papers comparing models fit through AIC (Akaike Information Criterion) and 4 papers mentioning the parametric Pearson residuals test. We conclude that, although the awareness of dispersion problems are steadily increasing in ecology, there is still the need for proper and more standardized tools for checking and testing them.

---

**Table 1.** Different types of dispersion evaluation and tests for GLMs and GLMMs with examples of available R packages and functions.

| Test | Principle | Details/Limitations | R package:: function | Supported models | References |
|---|---|---|---|---|---|
| **Likelihood Ratio Test (LRT)** | Compare two models with and without free dispersion parameters. Not a dispersion test. | Requires fitting two models, requires defining an alternative model. For example:<br>- Poisson and negative binomial or generalized Poisson<br>- binomial and beta-binomial | `pscl::odTest()` | GLM Poisson -> negative binomial with `MASS::glm.nb()` | Jackman (2024) |
| | | | `DCluster:: test.nb.pois()` | GLM Poisson -> negative binomial with `MASS::glm.nb()` | Lopez-Quílez (2005) |
| | | | `anova(…,test="LRT")` | Many GLM/GLMMs* | R Core Team (2024) |
| | | | `lmtest::lrtest()` | GLMs | Zeiles & Hothorn (2002) |
| **Score-like test** | Score test: Evaluate score of restricted dispersion parameter | Requires score calculation for specific models. R functions only for Poisson GLM. | `DCluster::DeanB()`<br>`DCluster::DeanB2()` | GLM Poisson.<br>Score tests based on Dean (1992) | Lopez-Quílez (2005) |
| | | | `Rfast2:: overdispreg.test()` | GLM Poisson (own model implementation) | Papadakis et al. (2025) |
| | Regression-based test for overdispersion from Cameron & Trivedi (1990) | Distribution specific (Poisson-based only). | `overdisp::overdisp()` | GLM Poisson (own model implementation) | Cameron & Trivedi (2023) |
| | | | `AER::dispersiontest()` | GLM Poisson from `stats::glm()` | Kleiber & Zeileis (2008) |
| **Standardized residuals dispersion** | A goodness-of-fit test to evaluate residual dispersion, e.g. via sum of Pearson residuals. | **Parametric Pearson residuals test:** Assume Pearson residuals are chi-squared distributed.<br>For complex models, difficult to define parametric null distribution (unclear residual degrees of freedom). | `msme::P__disp()` | GLMs | Hilbe & Robinson (2025) |
| | | | `aods3::gof()` | GLMs | Lesnoff et al. (2024) |
| | | | `DHARMa:: testDispersion(…, type="Pearson")` | GLMs/GLMMs (naïve residual *df*) | Hartig (2024) |
| | | | `performance:: check_overdispersion()` | GLMs/GLMMs (naïve residual *df*) | Lüdecke (2021) |
| | | | `RVAideMemoire:: overdisp.glmer()` | GLMMs (from `lme4` package, naïve residual *df*, calculates only dispersion statistic, no test) | Herve (2025) |
| | | **Nonparametric Pearson residuals test:** Parametric bootstrapping of the model to generate a nonparametric estimate of the null distribution of the Pearson statistic.Computational costly. | `DHARMa:: testDispersion(…, refit=T, type="Pearson")` | GLMs/GLMMs | Hartig (2024) |
| | | **Deviance residuals:** assumes the residual deviance are chi-squared distributed. | `aods3::gof()` | GLMs | Lesnoff et al. (2024) |
| | | **Dispersion metric:** the square root of the penalized residual sum of squares divided by the number of observations. | `blmeco:: dispersion_glmer()` | GLMMS from `lmer::glmer()` (computes dispersion parameter only, no test) | Korner-Nievergelt et al. (2019) |
| **Raw residual variance** | Compares the expected to the observed variance in raw residuals. | Expected variance of raw residuals calculated through simulations of fitted model.<br>Fast nonparametrics but possibly less exact than working on the standardized residual dispersion. | `DHARMa:: testDispersion(…, type="DHARMa")` | GLMs/GLMMs | Hartig (2024) |

\* Different packages have the S3 method for the *anova* function to perform LRT.

## A short review of existing approaches to dispersion tests

157

158       After reviewing the available literature, we divided the proposed strategies for

159 addressing dispersion problems into four classes (Table 1). Here, we discuss these broad

160 strategies in more detail and explain why we focused on two of these classes as the most

161 suitable competitors for a general dispersion test for GLMs and GLMMs. We note that,

162 in addition to the four approaches mentioned here, dispersion problems may also show

163 up in general goodness-of-fit tests (e.g., Feng et al., 2020). However, as they are not

164 specifically designed to react to dispersion, we did not consider them further.

165 *Likelihood ratio tests*

166       A first general strategy for detecting dispersion problems is to compare a model

167 with fixed dispersion to its nearest "relative" with variable dispersion using a likelihood

168 ratio test (LRT) or another model selection technique, such as AIC (Yang et al., 2007).

169 For count data, this could involve comparing a Poisson GLM to a negative binomial or

170 generalised Poisson GLM (Hilbe, 2014), or comparing a binomial GLM to a beta-

171 binomial GLM (Dunn & Smyth, 2018). While relatively easy to implement, the

172 downside of this approach, apart from the higher computational cost of fitting two

173 models, is that it doesn't provide any direct diagnostics of over- or underdispersion. The

174 alternative model, however, might also fit better or worse for reasons other than a

175 dispersion problem. Moreover, using LRTs to detect dispersion problems has also been

176 discouraged, as they may yield unreliable results (Dean, 1992) and tend to

177 underestimate the evidence against the base model (Lawless, 1987). Therefore, we do

178 not find this approach suitable as a general dispersion test and do not consider it further.

179 *Score tests*

180    A second traditional option for assessing overdispersion is the score test (Dean,

181    1992; Dean & Lawless, 1989; Lawless, 1987). Score tests, also known as Lagrange

182    Multiplier (LM) tests, evaluate the gradient of the log-likelihood (called the score or

183    LM statistic) of a restricted parameter estimator (e.g., an overdispersion estimator

184    constrained to zero). Under the null hypothesis that the overdispersion is indeed zero,

185    the score will have an asymptotic chi-squared distribution (Rao, 1948). In performance

186    comparisons, score tests have been found to have good power (Ohara Hines, 1997), but

187    their disadvantage is that they are usually model specific (in the sense that different tests

188    are needed for Poisson or binomial GLMs); their implementation can be

189    computationally demanding; and, as they require access to the score, they must usually

190    be implemented with the model and cannot be calculated on top of a fitted model object.

191    Perhaps because of these issues, we were unable to find any R function that computes

192    score tests beyond the Poisson GLM (Table 1), although score tests have been

193    developed for other models, such as the binomial GLM (Dean, 1992).

194    An equivalent test related to the score test under certain conditions is the

195    regression-based overdispersion test proposed by Cameron & Trivedi (1990). Under a

196    Poisson model, the squared deviation of the observations from their fitted mean, after

197    subtracting the observation itself and scaling by the fitted mean, has expectation zero. In

198    contrast, under the negative binomial, it increases systematically with the mean. This

199    motivates an auxiliary regression of the transformed variable against the fitted mean,

200    with a significant slope indicating extra-Poisson variation. The main advantage of this

201    test against other score tests is its ease of implementation: it can be carried out after

202    fitting a standard Poisson GLM. However, similar to an LRT, the linear regression

203    imposes a particular form of overdispersion as an alternative hypothesis, and therefore,

204    seems less general than the test based on Pearson residuals described below.

205     We discarded score tests in general, and the Cameron & Trivedi (1990) test in

206     particular, from our further analysis, as it seems impractical to implement them across a

207     wide range of existing GLMM software.

208     *Tests based on residual dispersion*

209     A third class of testing approaches, arguably the most intuitive, directly

210     calculates a test statistic or goodness-of-fit metric on standardised model residuals. The

211     most widely used test of this kind is based on the sum of the model's Pearson residuals.

212     As Pearson residuals divide the raw residuals by the expected residual standard

213     deviation, a correctly specified model is expected to have a Pearson residual of around 1

214     for each observation. A dispersion statistic is then defined as the sum of squared

215     Pearson residuals divided by the residual degrees of freedom. Models with a so-defined

216     dispersion statistic > 1 are considered overdispersed, while dispersion statistics < 1 are

217     underdispersed. Sometimes, this metric is modified by replacing the sum of squared

218     Pearson residuals with the model deviance, which is typically more readily available.

219     However, as Venables & Ripley (2002) discuss, this metric should be avoided, as it

220     often deviates from 1, even for correctly specified GLMs.

221     Defining dispersion via the Pearson statistic has the added advantage that for a

222     GLM, the expected distribution under the null hypothesis of a correctly specified model

223     asymptotically follows a chi-squared distribution (McCullagh, 1985). This allows a

224     straightforward construction of a hypothesis test, where we compare the Pearson

225     statistic to the chi-squared distribution with the respective residual degrees of freedom

226     (*df*). This test is referred to with different terminologies, such as the Pearson chi-squared

227     dispersion test, the Pearson residuals-based test for overdispersion, or simply the

228  Pearson dispersion test. Hereafter, we refer to this test as the **parametric Pearson**

229  **residuals test** to differentiate it from the nonparametric version discussed below.

230      An alternative approach to constructing a dispersion test based on the Pearson

231  dispersion statistic involves generating a null distribution through parametric

232  bootstrapping. A parametric bootstrap means that new data are simulated from the fitted

233  model, and then the statistic of interest (in this case: the Pearson statistic of a fitted

234  model) is calculated based on these data. The parametric bootstrap has been previously

235  used for hypothesis tests in mixed-effects models where parametric null distributions

236  were difficult to obtain (e.g., Barr et al., 2013; Luke, 2017), and thus it seems a logical

237  alternative for more complicated models where the chi-squared distribution of the

238  Pearson dispersion statistic cannot be taken for granted (see methods for GLMMs

239  below). Nevertheless, implementing parametric bootstrapping in complex models can

240  be less efficient for at least two reasons: it is time-consuming and prone to errors in

241  model refits (Luke, 2017; Moral et al., 2017). A dispersion test based on this principle

242  was implemented in R by Hartig (2024). Hereafter, we will refer to this test as the

243  **nonparametric Pearson residuals test**.

244  *Tests based on simulated residual variances*

245      Simulation approaches can also be useful for generating null distributions for

246  alternative dispersion metrics. A last class of dispersion test approaches, which, to our

247  knowledge, was introduced in the DHARMa R package (Hartig, 2024) but has not been

248  discussed in the literature so far, involves defining a test statistic based on the dispersion

249  of the raw residuals. The test compares the observed raw residual variance (differences

250  between the observed data and the model predictions) with the simulated raw residual

251  variances (differences between the simulated data and the model predictions). Both

252 variances are scaled to the variance of all simulated observations to account for

253 differences in the number of simulations across the fitted model. For GLMMs, data

254 simulations can be generated conditionally or unconditionally on the fitted random

255 effects. The dispersion statistic is then defined as the ratio of the observed residual

256 variance to the mean of the simulated residual variances. Similar to the Pearson statistic,

257 a ratio > 1 indicates overdispersion, a ratio < 1 indicates underdispersion, and a

258 significance test is constructed based on the distribution of simulated residual variances.

259     From a theoretical perspective , this approach seems less elegant than the use of

260 Pearson residuals, because the latter, by "standardising" the residual dispersion relative

261 to the expected dispersion, allows each data point to contribute similarly to the

262 dispersion statistic. In contrast, the test on the unstandardised residuals will be more

263 influenced by large data points. However, the primary advantage of this approach is

264 computational, as it enables a nonparametric estimate of the test statistic without

265 requiring a re-fit of the model (in contrast to the nonparametric Pearson residuals test).

266 Hereafter, we will refer to this test as the **simulation-based residual variance test** to

267 differentiate it from tests based on Pearson residuals.

268 **Methods**

269 *Selected models and setup of the performance comparisons*

270 After reviewing the available approaches, we identified three tests as potential

271 candidates for a generally applicable dispersion test that could be implemented across a

272 wide range of GLMs and GLMMs:

273     (1) The parametric Pearson residuals test

274     (2) The nonparametric Pearson residuals test

275     (3) The simulation-based residual variance test

276         To compare the performance of these three tests, we simulated datasets based on

277     the two main distributions that often exhibit over- or underdispersion: the Poisson and

278     the binomial (N/K) proportions. We varied the sample size (from 10 to 10,000) and the

279     intercept (from -3 to 3, at the link function scale) for the simulated data from both

280     distributions. We simulated a gradient of overdispersed data by adding noise to the

281     linear predictor with values from a Gaussian distribution with a mean of zero and ten

282     standard deviation values varying from 0 to 1. We evaluated test performance by

283     comparing type I error, power, and dispersion statistics across all parameter

284     combinations in the simulated datasets.

285         All models were fitted using the functions *glm* from the stats package or *glmer*

286     from the lme4 package (Bates et al., 2015) in R (v4.4; R Core Team, 2024). All

287     dispersion tests were performed with the DHARMa package (Hartig, 2024). For the

288     simulation-based residual variance test and the nonparametric Pearson residuals test, we

289     set the number of simulations to 250 (the default in DHARMa). All simulations and

290     analysis codes are available at this repository

291     (https://anonymous.4open.science/r/dispersion_test_GLMM/README.md). The

292     supplementary material provides a script file with instructions and examples for

293     applying dispersion tests using the DHARMa package.

294     *Theoretical expectations*

295         The classical (1) parametric Pearson residuals test assumes that the sample size

296     (n-asymptotic) and the expected values are sufficiently large (phi-asymptotic) (Venables

297     & Ripley, 2002). This implies that when the expected counts (or intercept) and/or the

298     number of observations are small, Pearson residuals may not provide reliable

299 information about model fit (see S2). Some corrections for Pearson residuals in small

300 samples have been suggested (e.g., Cordeiro, 2004; Cordeiro & Simas, 2009), but they

301 are not currently implemented in the most common R packages. Therefore, we expect

302 the parametric Pearson residuals test to perform well for GLMs, except in very small

303 sample sizes and expected counts (hereafter "small-data" situations).

304      It is unclear whether the parametric Pearson residuals test can be extended to

305 GLMMs or other hierarchical models, where counting residual degrees of freedom ($df$)

306 is not straightforward (Bolker et al., 2009; Luke, 2017). In mixed-effects models, the $df$

307 associated with a random effect are data-specific (adaptive shrinkage) and expected to

308 lie between one and the number of grouping factors (Baayen et al., 2008; Bolker et al.,

309 2009; Luke, 2017). Approaches exist to approximate $df$ for random effects in LMMs

310 (e.g., Schaalje et al., 2002), but their generalisation to GLMMs remains an active area

311 of research. Current R packages that implement the parametric Pearson residuals test

312 approximate the $df$ using the so-called naïve $df$ (e.g., n = 1 per random effect) for testing

313 LMMs/GLMMs (Table 1). We expect that the error introduced by this approximation

314 increases with the number of random-effect groups. To test this, we varied the number

315 of groups in the random intercept (10, 50, and 100) in our simulated data.

316      In contrast to the parametric Pearson residuals test, we expect the (2)

317 nonparametric Pearson residuals test to be robust to small-data problems and to the

318 presence of random effects, as it doesn't rely on a specific parametric distribution.

319 However, because the test uses parametric bootstrapping, we expected it to run much

320 more slowly than the other tests, especially for more complex GLMMs. For this

321 purpose, we compared the runtimes of the tests using a small set of simulated data (see

322 S7).

323       For GLMMs tested with the (3) simulation-based residual variance test, we

324    compared the test's performance under the two simulation approaches, conditional and

325    unconditional on random effects. We expect lower power in the unconditional

326    simulation results, as overdispersion is a phenomenon at the model distribution level

327    (i.e., at a higher level). We evaluated the circumstances under which this test is reliable

328    as a fast alternative to both dispersion tests based on Pearson residuals.

329    **Results**

330    *Performance on Poisson and binomial GLMs*

331       For Poisson GLMs, we found the expected distribution problems (Fig. 2): type I

332    error rates for the parametric Pearson residuals test were substantially high for the

333    smallest intercepts (-3), and they did not reach the nominal value of 0.05 even for very

334    large sample sizes (n = 10,000). The type I error rates for the nonparametric Pearson

335    residuals test were well calibrated, except for the smallest intercept (-3), with slightly

336    conservative type I error rates (< 0.05). For the simulation-based residual variance test,

337    type I errors were independent of sample size, but exhibited an intercept-dependent

338    conservative bias, ranging from almost 0 for the smallest intercept to 0.06 for the largest

339    intercept.

340       For binomial GLMs, the type I error rates for the parametric Pearson residuals

341    test were generally conservatively calibrated around 0.04 (Fig. 2). Type I error rates for

342    the nonparametric Pearson residuals test averaged around 0.05 and 0.06, except for the

343    very low and very high intercepts (-3 and 3). For the simulation-based residual variance

344    test, type I error rates were conservatively very low for all simulated parameters,

345    bouncing below 0.01.

**Figure 2.** The simulation-based residual variance test has a more conservative type I error rate than both Pearson residual tests. The three dispersion tests were applied to Poisson (upper panels) and binomial proportion (lower panels) GLMs: 1a) parametric Pearson residuals test, 1b) nonparametric Pearson residuals test, and 2) simulation-based residual variance test (see Table 1 for explanations). Simulations were run across different sample sizes (x-axis) and intercepts (colours, values on the link function scale). In B), the model is a binomial proportion with ten trials. All points include a 95% confidence interval calculated from exact binomial tests across the 10,000 simulations. Note the square-root scale of the y-axis in plot A. The dotted horizontal black line shows the 0.05 nominal type I error value.

The statistical power of the simulation-based residual variance test was lower than the parametric and nonparametric Pearson residuals tests for both binomial and Poisson GLMs, but tended to be similar with larger sample sizes (Fig. 3). We found that the reason for this is the very conservative type I error rates (Fig. 2). When power is calibrated by using the p-value at the 5% quantile of its empirical distribution for each simulation (details in S6), the differences disapear (Fig. 3).

The dispersion statistics of the simulation-based residual variance test were highly dependent on the intercept, slope, and number of trials in the binomial model (see S5, Fig. S5.2), and tended to be smaller than those based on Pearson residuals. In

366   contrast, for Poisson models, the values tended to be larger than those of Pearson

367   statistics (Fig. S4.5). This may also explain the lower uncorrected power for the

368   simulation-based residual variance test, especially for binomial models.



369

**Figure 3.** The simulation-based residual variance test (in yellow) has lower power than both Pearson residuals tests (green and blue) for GLMs unless power is calibrated by type I error rates (dashed lines). Lower power is more evident for binomial models (upper panel) and smaller sample sizes (first two columns). Results based on 10,000 simulations per combination of parameters for an intercept = 0 and slope = 1. For all simulation results, see Fig. S6.1 and S6.2.

376   *GLMM performance*

377       For the GLMMs, we first compared the performance of the parametric Pearson

378   residuals test (two-sided) for an increasing number of groups ($m$) in the random

379   intercepts. As expected, the performance of the test failed for a large number of groups

380   in the random effects (Fig. 4A). The dispersion statistic was underestimated, and the

381   type I error rates were too high because the test detected significant underdispersion.

382   Testing only for overdispersion ("greater" test) when using the parametric Pearson

383  residuals test appears to be the only reasonable approach for GLMMs (Fig. 4A). Still, it

384  doesn't prevent the dispersion statistics from being biased to lower values.



**Pearson Chi-squared dispersion tests for GLMMs**

385

386  **Figure 4**. The parametric Pearson residuals test failed for GLMMs with many groups in
387  the random intercepts (plot panels). A) Power and type I error rates (blue shaded area)
388  for the "two-sided" (solid lines) and "greater" (dotted lines) chi-squared tests for the
389  Pearson statistic. B) Pearson dispersion statistics with the red shaded area indicating
390  dispersion statistics estimated below 1 (underdispersion). Notice that the y-axis of plot
391  B is on a logarithmic scale of 10. Results with 10,000 simulations for an intercept of 0
392  and a sample size (n) of 1,000 data points.

393       When comparing the alternative dispersion tests for GLMMs, the nonparametric

394  Pearson residuals test presented very good results, with a type I error rate around 0.05

395  (Fig. S7.1 and S7.2) and higher power than the simulation-based residual variance tests

396  (Fig. 5). As expected, the unconditional simulation-based residual variance test had the

397  worst performance: very low type I errors (Fig. S7.1 and S7.2), very low power, and

398  dispersion statistics below 1 (Fig. 5B), especially for Poisson models. The conditional

399  simulation-based residual variance test also had very small type I errors (Fig. S7.1 and

400  S7.2), but power increased with the simulated overdispersion. The performance of both

401 simulation-based residual variance tests (unconditional and conditional) didn't change

402 much with the number of groups for the Poisson GLMMs, but it improved for the

403 binomial GLMMs with the increasing number of groups in the random intercept.



**Figure 5**. The nonparametric Pearson residuals test showed correct Type 1 error, higher
power, and larger dispersion statistics than the simulated-based residual variance tests
(conditional and unconditional to all random effects) for Poisson and binomial GLMMs.
Power (A), type I error (shaded blue area in A), and dispersion statistics (B) for the
alternative dispersion tests for Poisson and binomial GLMMs with different numbers of
groups in random intercepts. The dashed horizontal line in (A) indicates the nominal

411 value of 0.05 for type I error. The dotted horizontal line in (A) indicates the 50% power,
412 and the dotted horizontal line in (B) indicates the dispersion statistics of 1. The results
413 are based on 1,000 simulations per parameter combination, with an intercept of 0 and a
414 sample size (n) of 1,000.

## Discussion

416      The goal of this study was to identify a dispersion test that is widely applicable

417 across different GLM and GLMM distributions and random-effects structures, which

418 are common in ecological data analysis. Our conclusion is that the nonparametric

419 Pearson residuals test is the most reliable general test currently available. For GLMs,

420 this test exhibited similar power to the parametric Pearson residuals test but with more

421 reliable type I error rates in small-sample situations. The downside of this test is that it

422 can be computationally expensive, with runtimes in the order of minutes for larger

423 GLMMs.

424      The simulation-based residual variance test for GLMs is fast to compute, but its

425 dispersion statistic is more difficult to interpret and often leads to overly conservative

426 type I errors. This results in low power unless it is additionally calibrated using a

427 simulated p-value distribution. The parametric Pearson residuals test is computationally

428 efficient, but it is unreliable in small-data situations and in the presence of random

429 effects. Below, we discuss these points in more detail and provide recommendations for

430 general users who rely on already implemented R packages for model fit and

431 diagnostics.

432 *Why and when does the parametric Pearson residuals test fail?*

433      We showed that the parametric Pearson residuals test, although popular, quick,

434 and relatively easy to compute, has two main disadvantages: it performs poorly in (1)

435 small-data situations (Fig. 2) and (2) in the presence of random effects (Fig. 4). The first

436 problem arises from a mismatch between the distribution of the Pearson statistic and the

437 chi-squared distribution under small-data conditions (Fig. S2.1 and S2.2). This

438 phenomenon has already been studied (e.g., Fletcher, 2012; Kuss, 2002), with suggested

439 corrections (Farrington, 1996; McCullagh, 1985). However, none of these corrections

440 are implemented in the current R packages (Table 1), and we believe it will be difficult

441 to devise corrections that work across a wide range of distributions.

442 The second problem arises because counting 1 degree of freedom ($df$) for a

443 random effect, as done in most implementations of this test, typically underestimates the

444 true model $df$, and this underestimation increases with the random effect's number of

445 levels. The result is a bias in the dispersion statistic towards underdispersion that

446 increases with the number of random-effect levels (Fig. 4). Two-sided tests would

447 therefore often wrongly detect significant underdispersion in perfectly valid GLMMs,

448 which is likely why most R implementations of this test only test for overdispersion.

449 When applying this test to GLMMs, we recommend following the same approach and

450 ignoring dispersion statistics smaller than 1. Nevertheless, this is an unsatisfactory

451 solution, as the biased dispersion statistic also causes a loss of power.

452 A possible solution for GLMMs could be to use a better approximation of the

453 residual degrees of freedom ($df$). For LMMs, approximations for denominator $df$ have

454 been successfully used for hypothesis testing (Luke, 2017), for example, the

455 Satterthwaite (1946) and the Kenward-Roger (2009). Although there is some evidence

456 that these approximations are also accurate for GLMMs (Stroup, 2015), the main R

457 packages implementing these methods are currently limited to LMMs (e.g., *pbkrtest*

458 Halekoh & Højsgaard, 2014; *lmerTest* Kuznetsova et al., 2017). However, the recently

459 released package *glmmrBase* (Watson, 2024) allows these methods to be applied to

460 GLMMs. We performed some parametric Pearson residuals tests for Poisson GLMMs

461     using a modified residual *df* approximation (see S9). Although the parametric Pearson

462     residuals tests with the approximated residual *df* performed much better than those with

463     the naïve residual *df*, they still underperformed compared to the nonparametric Pearson

464     residuals test when there were a large number of groups in the random effects (Fig.

465     S9.4), especially in small-data situations.

466     *When are simulation-based residual variance tests an alternative?*

467           The simulation-based residual variance test developed in the R package

468     DHARMa (Hartig, 2024) is the main alternative to the family of Pearson residuals tests.

469     Its principle is simple: when the model is correctly specified, the variance of the

470     observed data should match that of data simulated from the model. The main advantage

471     of this approach is that it is a non-parametric test applicable to any model structure and

472     does not require refitting the model, making it considerably faster and easier to

473     implement in statistical software. We also note that for GLMMs, simulations should be

474     performed conditionally to avoid a loss of power, presumably due to the increased

475     variability created by re-simulating the random effects (unconditional simulations).

476           The disadvantages of this approach are that it is often overly conservative,

477     resulting in lower power than the Pearson residuals tests. Additionally, the calculated

478     dispersion statistic differs from the Pearson dispersion statistic, making it difficult to

479     compare the two approaches. We conjectured that both problems could be related to the

480     test statistic being based on the raw variance (rather than a scaled variance, as with the

481     Pearson statistics), which may overrepresent observations with large values. We

482     considered scaling each observation by the expected variance, but this is not readily

483     available for a wide class of models, and using simulations to approximate it fails for

484     discrete-valued distributions (see S8).

## Conclusions and recommendations

486        Although neither of the considered options for testing dispersion excelled in all

487    dimensions (Fig. 6), our primary recommendation is that for standard GLMs with

488    sufficient data, the parametric Pearson chi-squared test, available in many packages

489    (Table 1), can be safely used. In complex situations, particularly for GLMMs, we

490    recommend the nonparametric Pearson residuals test. It has very few weaknesses, other

491    than being computationally costly. If the nonparametric Pearson residuals test cannot be

492    calculated due to speed or convergence problems with refitting complex models, we

493    recommend using the simulation-based residual variance test with simulations

494    performed conditionally on the fitted random effects. All three approaches are available

495    via the *testDispersion* function in the DHARMa R package (Hartig, 2024). We provide

496    a supplementary file with instructions and an example for applying dispersion tests

497    using the DHARMa package.



| | GLM | GLM ("small-data") | GLMM (few RE groups) | GLMM (many RE groups) | Speed |
|---|---|---|---|---|---|
| Simulation-based residual variance | ++ | - | ++ | + | + |
| Nonparametric Pearson residuals | ++ | + | ++ | ++ | - |
| Parametric Pearson residuals | ++ | - | + | - | ++ |

498

**Figure 6.** Performance comparisons of the dispersion tests evaluated for each
"dimension" for Poisson and binomial models: GLMs in general, GLMs with small
sample size or intercept ("small data"), GLMMs with one random effect with few
groups/levels, GLMMs with many groups/levels in a random effect, and computational
time for calculating the test (speed). The symbols mean: "-" bad performance, "+" good
performance, "++" very good performance.

505      Although our simulation examples focused on overdispersion, the tests

506      considered in our study can also be used to detect underdispersion by testing the

507      dispersion "two-sided" or "less than" against null statistics. The clear exception would

508      be testing for underdispersion using the parametric Pearson residuals test for GLMMs,

509      which would be anti-conservative due to the discussed bias towards underdispersion in

510      the presence of random effects.

511      *Recommendations for ecological data analysis when using dispersion tests*

512      For interpretation and applied ecological data analysis, we stress that a

513      significant over- or underdispersion result does not necessarily indicate that the

514      distribution must be changed. First, hypothesis tests evaluate statistical rather than

515      ecological significance. In other words, a significant test for overdispersion indicates

516      that the overdispersion signal deviates from a null expectation, but the p-value does not

517      measure the strength of the deviation. The first step in a dispersion test should thus be to

518      examine how much the dispersion statistic deviates from the expected value of 1. For

519      very large sample sizes, small departures from 1 may be statistically significant, but

520      they may not necessarily warrant a change to the model. Second, after finding that a

521      dispersion problem is both significant and meaningful, we suggest checking for

522      problems beyond the distribution, such as heteroscedasticity, missing predictors, an

523      incorrect link function, excess zeros, or overfitting. In our experience, these types of

524      model misspecifications often cause over-/underdispersion, but can be distinguished

525      from a "real" distributional problem through careful residual checks. Blindly changing

526      the distribution only masks the problem, without offering a real solution to the

527      underlying problems.

528     Finally, after ruling out potential model misspecifications leading to under-

529     /overdispersion, we should consider changing the model's distribution, as we may be

530     facing an 'intrinsic' under-/overdispersion problem, likely due to the nature of

531     ecological data (Box 1). A traditional and flexible solution is to use the 'quasi'

532     distributions (Wedderburn, 1974), which essentially correct p-values but do not

533     represent an explicit data-generating process with an associated likelihood, precluding,

534     for example, simulation from the fitted model. A second alternative for adding

535     dispersion is to use observation-level random effects (Bolker et al., 2009; Elston et al.,

536     2001; Harrison, 2014; Ozgul et al., 2009). While often a reasonable solution, excessive

537     use of random effects can create problems in calculating other statistical indicators

538     (such as p-values) that we would rather avoid. For that reason, we consider the best

539     solution to address 'intrinsic' under-/overdispersion is to switch to the corresponding

540     variable-dispersion distributions. For overdispersed count data, the most used is the

541     negative binomial (see S1). However, other distributions have been used in ecology to

542     handle both over- and underdispersion, such as the generalised Poisson, the Conway-

543     Maxwell-Poisson, the Double Poisson, and the Good distributions (Agis et al., 2024;

544     Brooks et al., 2019; Lynch et al., 2014). For discrete proportions data, the beta-binomial

545     distribution (Harrison, 2015) is considered the most appropriate for overdispersed

546     binomial models (Harrison, 2015). Regardless of the approach, an "over-

547     /underdispersion-free" GLM/GLMM is essential for better interpretation of ecological

548     models and for facilitating sound scientific discoveries.


549     **References**

550     Agis, D., Tur, J., Moriña, D., Puig, P., & Fernández-Fontelo, A. (2024). good: An R
551            package for modelling count data. *Methods in Ecology and Evolution*, *15*(12),
552            2192–2197. https://doi.org/10.1111/2041-210X.14387

Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language, Special Issue: Emerging Data Analysis*, *59*(4), 390–412. https://doi.org/10.1016/j.jml.2007.12.005

Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, *68*(3), 255–278. https://doi.org/10.1016/j.jml.2012.11.001

Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software*, *67*(1), 1–48. https://doi.org/10.18637/jss.v067.i01

Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H. H., & White, J.-S. S. (2009). Generalized linear mixed models: A practical guide for ecology and evolution. *Trends in Ecology & Evolution*, *24*(3), 127–135. https://doi.org/10.1016/j.tree.2008.10.008

Brooks, M. E., Kristensen, K., Darrigo, M. R., Rubim, P., Uriarte, M., Bruna, E., & Bolker, B. M. (2019). Statistical modeling of patterns in annual reproductive rates. *Ecology*, *100*(7), e02706. https://doi.org/10.1002/ecy.2706

Cameron, A. C., & Trivedi, P. (2023). *overdisp: Overdispersion in count data multiple regression analysis* (Version 0.1.2) [Computer software]. https://CRAN.R-project.org/package=overdisp

Cameron, A. C., & Trivedi, P. K. (1990). Regression-based tests for overdispersion in the Poisson model. *Journal of Econometrics*, *46*(3), 347–364. https://doi.org/10.1016/0304-4076(90)90014-K

Campbell, H. (2021). The consequences of checking for zero-inflation and overdispersion in the analysis of count data. *Methods in Ecology and Evolution*, *12*(4), 665–680. https://doi.org/10.1111/2041-210X.13559

Campbell, N. A., & Reece, J. B. (2005). *Biology* (7th ed.). Pearson Benjamin Cummings.

Collings, B. J., & Margolin, B. H. (1985). Testing Goodness of Fit for the Poisson Assumption When Observations Are Not Identically Distributed. *Journal of the American Statistical Association*, *80*(390), 411–418.

Cordeiro, G. M. (2004). On Pearson's residuals in generalized linear models. *Statistics & Probability Letters*, *66*(3), 213–219. https://doi.org/10.1016/j.spl.2003.09.004

Cordeiro, G. M., & Simas, A. B. (2009). The distribution of Pearson residuals in generalized linear models. *Computational Statistics & Data Analysis*, *53*(9), 3397–3411. https://doi.org/10.1016/j.csda.2009.02.025

Dean. (1992). Testing for Overdispersion in Poisson and Binomial Regression Models. *Journal of the American Statistical Association*, *87*(418), 451–457. https://doi.org/10.2307/2290276

Dean, C., & Lawless, J. F. (1989). Tests for Detecting Overdispersion in Poisson Regression Models. *Journal of the American Statistical Association*, *84*(406), 467–472. https://doi.org/10.1080/01621459.1989.10478792

Dunn, P. K., & Smyth, G. K. (2018). *Generalized Linear Models With Examples in R*. Springer New York. https://doi.org/10.1007/978-1-4419-0118-7

Elston, D. A., Moss, R., Boulinier, T., Arrowsmith, C., & Lambin, X. (2001). Analysis of aggregation, a worked example: Numbers of ticks on red grouse chicks. *Parasitology*, *122*(05), 563–569.

Farrington, C. P. (1996). On Assessing Goodness of Fit of Generalized Linear Models to Sparse Data. *Journal of the Royal Statistical Society. Series B (Methodological)*, *58*(2), 349–360.

Feng, C., Li, L., & Sadeghpour, A. (2020). A comparison of residual diagnosis tools for diagnosing regression models for count data. *BMC Medical Research Methodology*, *20*(1), 175. https://doi.org/10.1186/s12874-020-01055-2

Fisher, R. A. (1950). The Significance of Deviations from Expectation in a Poisson Series. *Biometrics*, *6*(1), 17–24. https://doi.org/10.2307/3001420

Fletcher, D. J. (2012). Estimating overdispersion when fitting a generalized linear model to sparse data. *Biometrika*, *99*(1), 230–237. https://doi.org/10.1093/biomet/asr083

Halekoh, U., & Højsgaard, S. (2014). A Kenward-Roger Approximation and Parametric Bootstrap Methods for Tests in Linear Mixed Models – The R Package pbkrtest. *Journal of Statistical Software*, *59*, 1–32. https://doi.org/10.18637/jss.v059.i09

Harrison, X. A. (2014). Using observation-level random effects to model overdispersion in count data in ecology and evolution. *PeerJ*, *2*, e616. https://doi.org/10.7717/peerj.616

Harrison, X. A. (2015). A comparison of observation-level random effect and Beta-Binomial models for modelling overdispersion in Binomial data in ecology & evolution. *PeerJ*, *3*, e1114. https://doi.org/10.7717/peerj.1114

Hartig, F. (2024). *DHARMa: Residual Diagnostics for Hierarchical (Multi-Level / Mixed) Regression Models* (Version 0.4.7) [Computer software]. https://CRAN.R-project.org/package=DHARMa

Herve, M. (2025). *RVAideMemoire: Testing and plotting procedures for biostatistics* (Version 0.9-83-11) [Computer software]. https://CRAN.R-project.org/package=RVAideMemoire

Hilbe, J. M. (2011). *Negative Binomial Regression*.

Hilbe, J. M. (2014). *Modeling count data*. Cambridge University Press.

Hilbe, J., & Robinson, A. (2025). *msme: Functions and datasets for "methods of statistical model estimation"* (Version 0.5.4) [Computer software]. https://CRAN.R-project.org/package=msme

Jackman, S. (2024). *pscl: Classes and methods for R developed in the political science computational laboratory* (Version 1.5.9) [Computer software]. University of Sydney. https://github.com/atahk/pscl/

Kenward, M. G., & Roger, J. H. (2009). An improved approximation to the precision of fixed effects from restricted maximum likelihood. *Computational Statistics & Data Analysis*, *53*(7), 2583–2595. https://doi.org/10.1016/j.csda.2008.12.013

Kleiber, C., & Zeileis, A. (2008). *Applied econometrics with R [R package AER version 1.2-14]*. Springer-Verlag. https://doi.org/10.1007/978-0-387-77318-6

639 Korner-Nievergelt, F., Roth, T., Felten, S. von, Guelat, J., Almasi, B., & Korner-
640     Nievergelt, P. (2019). *blmeco: Data Files and Functions Accompanying the*
641     *Book "Bayesian Data Analysis in Ecology using R, BUGS and Stan"* (Version
642     1.4) [Computer software]. https://cran.r-
643     project.org/web/packages/blmeco/index.html

644 Kuss, O. (2002). Global goodness-of-fit tests in logistic regression with sparse data.
645     *Statistics in Medicine*, *21*(24), 3789–3801. https://doi.org/10.1002/sim.1421

646 Kuznetsova, A., Brockhoff, P., & Christensen, R. (2017). lmerTest Package: Tests in
647     Linear Mixed Effects Models. *Journal of Statistical Software, Articles*, *82*(13).
648     https://doi.org/10.18637/JSS.V082.I13

649 Lai, J., Lortie, C. J., Muenchen, R. A., Yang, J., & Ma, K. (2019). Evaluating the
650     popularity of R in ecology. *Ecosphere*, *10*(1), e02567.
651     https://doi.org/10.1002/ecs2.2567

652 Lawless, J. F. (1987). Negative binomial and mixed Poisson regression. *Canadian*
653     *Journal of Statistics*, *15*(3), 209–225. https://doi.org/10.2307/3314912

654 Lesnoff, M., Lancelot, R., & Siberchicot, A. (2024). *aods3: Analysis of Overdispersed*
655     *Data using S3 Methods* (Version 0.5) [Computer software]. https://cran.r-
656     project.org/web/packages/aods3/index.html

657 Lindén, A., & Mäntyniemi, S. (2011). Using the negative binomial distribution to model
658     overdispersion in ecological count data. *Ecology*, *92*(7), 1414–1421.
659     https://doi.org/10.1890/10-1831.1

660 Lopez-Quílez, V. G.-R. J. F.-F. A. (2005). Detecting clusters of disease with R. *Journal*
661     *of Geographical Systems*, *7*(2), 189–206. https://doi.org/10.1007/s10109-005-
662     0156-5

663 Lüdecke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021).
664     performance: An R package for assessment, comparison and testing of statistical
665     models. *Journal of Open Source Software*, *6*(60), 3139.
666     https://doi.org/10.21105/joss.03139

667 Luke, S. G. (2017). Evaluating significance in linear mixed-effects models in R.
668     *Behavior Research Methods*, *49*(4), 1494–1502. https://doi.org/10.3758/s13428-
669     016-0809-y

670 Lynch, H. J., Thorson, J. T., & Shelton, A. O. (2014). Dealing with under- and over-
671     dispersed count data in life history, spatial, and community ecology. *Ecology*,
672     *95*(11), 3173–3180. https://doi.org/10.1890/13-1912.1

673 McCullagh, P. (1985). On the Asymptotic Distribution of Pearson's Statistic in Linear
674     Exponential-Family Models. *International Statistical Review / Revue*
675     *Internationale de Statistique*, *53*(1), 61–67. https://doi.org/10.2307/1402880

676 McCullagh, P., & Nelder, J. (1989). Generalized linear models. *Journal of the Royal*
677     *Statistical Society*, *135*(3), 370–384.

678 McMahon, S. M., & Diez, J. M. (2007). Scales of association: Hierarchical linear
679     models and the measurement of ecological systems. *Ecology Letters*, *10*(6),
680     437–452. https://doi.org/10.1111/j.1461-0248.2007.01036.x

Moral, R. A., Hinde, J., & Demétrio, C. G. B. (2017). Half-Normal Plots and Overdispersed Models in R: The hnp Package. *Journal of Statistical Software*, *81*, 1–23. https://doi.org/10.18637/jss.v081.i10

Nakagawa, S., Ortega, S., Gazzea, E., Lagisz, M., Lenz, A., Lundgren, E., & Mizuno, A. (2026). Location–scale models in ecology and evolution: Heteroscedasticity in continuous, count and proportion data. *Methods in Ecology and Evolution*, *17*(2), 554–566. https://doi.org/10.1111/2041-210x.70203

Ohara Hines, R. J. (1997). A comparison of tests for overdispersion in generalized linear models. *Journal of Statistical Computation and Simulation*, *58*(4), 323–342. https://doi.org/10.1080/00949659708811838

Ozgul, A., Oli, M. K., Bolker, B. M., & Perez-Heydrich, C. (2009). Upper respiratory tract disease, force of infection, and effects on survival of gopher tortoises. *Ecological Applications*, *19*(3), 786–798.

Papadakis, M., Tsagris, M., Fafalios, S., Dimitriadis, M., & Lasithiotakis, M. (2025). *Rfast2: A collection of efficient and extremely fast R functions II* (Version 0.1.5.4) [Computer software]. https://CRAN.R-project.org/package=Rfast2

Puig, P., Valero, J., & Fernández-Fontelo, A. (2024). Some mechanisms leading to underdispersion: Old and new proposals. *Scandinavian Journal of Statistics*, *51*(1), 245–267. https://doi.org/10.1111/sjos.12677

Quine, M. P., & Seneta, E. (1987). Bortkiewicz's Data and the Law of Small Numbers. *International Statistical Review / Revue Internationale de Statistique*, *55*(2), 173–181. https://doi.org/10.2307/1403193

R Core Team. (2024). *R: a language and environment for statistical computing* (Version v4.4.1) [Computer software]. R Foundation for Statistical Computing. https://www.R-project.org/

Rao, C. R. (1948). Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation. *Mathematical Proceedings of the Cambridge Philosophical Society*, *44*(1), 50–57. https://doi.org/10.1017/S0305004100023987

Rhodes, J. R. (2015). Mixture models for overdispersed data. In G. A. Fox, V. J. Sosa, & S. M. Negrete-Yankelevich, *Ecological Statistics: Contemporary theory and applications*. Oxford University Press.

Satterthwaite, F. E. (1946). An Approximate Distribution of Estimates of Variance Components. *Biometrics Bulletin*, *2*(6), 110–114. https://doi.org/10.2307/3002019

Schaalje, G. B., McBride, J. B., & Fellingham, G. W. (2002). Adequacy of approximations to distributions of test statistics in complex mixed linear models. *Journal of Agricultural, Biological, and Environmental Statistics*, *7*(4), 512–524. https://doi.org/10.1198/108571102726

Stroup, W. W. (2015). Rethinking the Analysis of Non-Normal Data in Plant and Soil Science. *Agronomy Journal*, *107*(2), 811–827. https://doi.org/10.2134/agronj2013.0342

Touchon, J. C., & McCoy, M. W. (2016). The mismatch between current statistical practice and doctoral training in ecology. *Ecosphere*, *7*(8), e01394. https://doi.org/10.1002/ecs2.1394

726 Venables, W. N., & Ripley, B. D. (2002). *Modern Applied Statistics with S*. Springer
727     New York. https://doi.org/10.1007/978-0-387-21706-2

728 Watson, S. I. (2024). *Generalised Linear Mixed Model Specification, Analysis, Fitting,*
729     *and Optimal Design in R with the glmmr Packages* (arXiv:2303.12657). arXiv.
730     https://doi.org/10.48550/arXiv.2303.12657

731 Wedderburn, R. W. M. (1974). Quasi-likelihood functions, generalized linear models,
732     and the Gauss—Newton method. *Biometrika*, *61*(3), 439–447.
733     https://doi.org/10.1093/biomet/61.3.439

734 Xekalaki, E. (2014). On the distribution theory of over-dispersion. *Journal of Statistical*
735     *Distributions and Applications*, *1*(1), 19. https://doi.org/10.1186/s40488-014-
736     0019-z

737 Yang, Z., Hardin, J. W., Addy, C. L., & Vuong, Q. H. (2007). Testing Approaches for
738     Overdispersion in Poisson Regression versus the Generalized Poisson Model.
739     *Biometrical Journal*, *49*(4), 565–584. https://doi.org/10.1002/bimj.200610340

740 Zeileis, A., & Hothorn, T. (2002). Diagnostic checking in regression relationships. *R*
741     *News*, *2*(3), 7–10.

742

**Supporting information for:**

**Dispersion tests in generalised linear mixed-effects models - a**

**methods comparison and practical guide for ecologists**

## S1. Trend analysis and current ecological literature practices on dispersal issues

To understand the extent of ecological studies relying on GLMs/GLMMs for count and discrete proportion data and those that address dispersion issues, we conducted a text analysis of the ecological literature over the past 20 years. We used the R package 'europepmc' (v0.4.3, Jahn, 2023) to search for articles in the PubMed and Medline NLM databases from 2005 to 2025. We used combinations of words (Table S1.1) to retrieved the annual records for: (1) the percentage of ecological papers using GLMs/GLMMs for count and discrete proportion data (Figure S1.1a), (2) the percentage of those papers that mention dispersion terms in general (Figure S1.1b), (3) the percentage of ecological papers using GLMs/GLMMs for count data mentioning dispersion terms (Figure BOX 1, main text), and (4) the percentage of ecological papers using GLMs/GLMMs for discrete proportion that mentioning dispersion terms (Figure BOX 1, main text).

18 **Table S1.1.** Word combinations used for the literature review on ecological practices for
19 count and discrete proportion data analysed with GLMs/GMMs and dispersion issues.

| Terms | Words combination |
|---|---|
| 1. Ecology: | |
| | "ecology" OR "ecolog*" |
| 2. Generalised linear models for count and discrete proportion data: | |
| | "count data" OR "poisson" OR "negative binomial" OR "generalized poisson" OR "generalised poisson" OR "conway-maxwell poisson" OR "binomial" OR "beta-binomial" OR "binomial proportion" |
| 3. Generalized linear models for count data only: | |
| | "count data" OR "poisson" OR "negative binomial" OR "generalized poisson" OR "generalised poisson" OR "conway-maxwell poisson" |
| 4. Generalized linear models for discrete proportion data only | |
| | "binomial" OR "beta-binomial" OR "binomial proportion" |
| 5. Dispersion terms: | |
| | "overdispersion" OR "over dispersion" OR "over-dispersion" OR "underdispersion" OR "under dispersion" OR "under-dispersion" OR "dispersion" |

20

21       The percentage of papers that mention count or proportion data in the context of

22 GLM/GLMM analysis increased 4-fold over 20 years, but appears to have stabilised

23 since 2015 (Figure S1.1A). For those papers, there is an increasing trend in mentioning

24 dispersion terms, reaching almost 25% in 2025 (Figure S1.1B). However, it means that

25 3/4 of ecological papers that mentioned GLMs/GLMMs for analysing count and/or

26 discrete proportion still don't report checking for dispersion problems.

**A) GLMs/GLMMs within Ecology**

**B) Dispersion within GLMs/GLMMs**

27

**Figure S1.1.** Trend analysis from the last 20 years of (A) ecological papers mentioning GLM/GLMMs for count and/or discrete proportion data, and (B) ecological papers that use GLM/GLMMs and mention dispersion terms.

We then summarised the current practices in dispersal issues for the ecological studies using GLMs/GLMMs for count and discrete proportion by searching for papers with the combination of words of the groups 1, 2 and 5 (Table S1.1). The query retrieved 7634 articles; we further selected only open-access articles from journals that publish ecological papers in 2025. From the subset of 457 articles, we randomly selected 200 papers for detailed information searches and retrieved the first 100 papers within the scope (ecology) that used count or discrete proportion data. To reach 100 papers, we read 155 papers; 33 were out of scope, and 22 did not use count or discrete proportion data analysis. Among them, 89 papers explicitly mentioned a dispersion issue in the methods section; 81 papers mentioned overdispersion, 4 mentioned underdispersion, and 4 mentioned both (tested for both issues). A total of 69 papers explicitly reported checking for dispersion, whereas only 40 reported testing for dispersion problems or comparing model fits using AIC.

Of the 40 papers explicitly testing dispersion, 25 reported using the DHARMa R package (Hartig, 2024), and 5 reported using the performance package (Lüdecke et al.,

46    2021). However, almost all of them didn't mention which test. Model comparison using

47    AIC was reported in 5 papers, and the Pearson Chi-squared test (Pearson parametric

48    residuals test) in 4, including one paper that used GLMMs and reported underdispersion

49    in many models (Laumer et al., 2025). This recent literature review shows an increasing

50    number of ecological studies examining dispersion problems, underscoring the

51    importance of appropriate tools for their detection and testing.

52         Additionally, we found that the most common approach to address dispersion

53    issues in count data was to switch from the Poisson distribution to the negative

54    binomial, or starting with the negative binomial in the first place (46 out of 78 records,

55    59%). Only 3 papers used the generalised Poisson distribution, and 1 paper reported

56    using the Conway-Maxwell-Poisson for underdispersed data. The quasi-Poisson

57    approach was reported in 7 papers (9%), the use of an observation-level random effects

58    in a Poisson GLMM was reported in 5 papers (6%), and the use of a zero-inflated

59    (Poisson or negative binomial) model was reported 12 times (15%).

60         For discrete proportion data, we identified 7 papers that report alternative

61    modelling to account for overdispersion. The quasi-binomial approach and the beta-

62    binomial distribution were reported 3 times each. The use of an observation-level

63    random effects in a binomial GLMM was reported in just one paper.

## S2. Pearson statistics and Chi-squared distribution

For GLMs, the parametric Pearson residuals test assumes that the sample size (n-asymptotic) and the expected values are sufficiently large (phi-asymptotic). Therefore, when the expected counts (or intercept) and/or the number of observations are small, Pearson residuals may not provide reliable information about model fit. To test boundaries where Pearson statistics fail, we simulated data with very different sample sizes (from 10 to 10,000, depending on the simulation) and intercepts (from -3 to 3, at the link function scale) for Poisson and binomial proportion GLMs. For each distribution and parameter combination, we used the Kolmogorov-Smirnov test (KS test) of adherence to compare the empirical distribution of 1000 simulations of the Pearson residuals with the Chi-squared distribution having the same residual degrees of freedom. We repeated this procedure 100 times and recorded the proportion of significant KS tests.

For the Poisson GLMs, the Pearson statistics distribution clearly departed from the Chi-square distribution for very small intercepts (-3, -1.5) and sample sizes (10, 20 and 50) (Figure S2.1 A). Even for very large sample sizes (10,000), the distribution did not approximate the Chi-squared distribution for the smallest simulated intercept (-3). Consequently, the KS tests showed all significant results for all simulations with the intercept at -3, except for the largest sample size (10,000), where it decreased to 60%. As expected, the proportion of significant results decreased with sample size for intercepts at -1.5 and 0. For larger intercepts, it remained around 5% for all sample sizes (Figure S2.2A).

For the **binomial GLMs**, the Pearson statistics distribution clearly departed from the Chi-squared distribution for very small and large intercepts (-3, 3) and small

88    sample sizes (10, 20, 50) (Figure S2.1B). The proportion of significant KS tests

89    decreased with sample size, but did not reach the nominal value of 0.05, even for very

90    large sample sizes and intermediate intercept values (-1.5, 0, 1.5).



91

92    **Figure S2.1.**Proportion of significant Kolmogorov-Smirnov adherence tests between
93    the empirical distribution of 1000 simulations of the Pearson statistics and a Chi-
94    squared distribution with the same residual degrees of freedom for A) Poisson and B)
95    binomial GLMs. Proportions were calculated from 100 simulations for each
96    combination of the data parameters (sample size and intercept). For binomial data, the
97    number of trials was fixed at 10. The 95% confidence intervals (vertical lines) were
98    drawn from binomial exact tests for each result with $p = 0.05$.

**Pearson Statistics X Chi-squared distribution**

**Figure S2.2**. Mean Pearson statistics distribution (from 100 simulated curves) for the binomial (green) and Poisson (purple), and the Chi-square distribution in black.

## S3. Type I error rates for the GLMs

Figures S3.1 and S3.2 show the distribution of the p-values for the dispersion tests applied to the Poisson and binomial GLMs, respectively, with 10,000 simulations for each combination of intercept and sample size. For the dispersion tests with correct type I error rates around the nominal value of 0.05, the distributions of p-values should present a uniform distribution with density 1.

For the Poisson GLMs (Figure S3.1), the simulation-based residual variance test (in red) presented the largest departure of the expected distribution for the smallest intercepts (-3, -1.5) across all sample sizes. This explains why the type I error rates for the simulation-based residual tests were so low and varied according to the intercept but didn't change with the sample size (main text Figure 2A). The parametric Pearson test had the opposite pattern with very low p-values for the smallest intercept (-3), but it tended to approximate the uniform distribution (decreasing the peak for the low p-values) with sample size. The p-values for the nonparametric Pearson test also showed a departure from the uniform distribution for the smallest intercept (-3), but tended to approach the uniform distribution with larger sample sizes and intercepts.

For the binomial GLMs (Figures S3.2), the p-values distribution of the simulation-based residual variance test also presented the largest departure from the uniform distribution, but for all intercepts and sample sizes. The p-values for both parametric Pearson and nonparametric Pearson tests were similar and tended towards the uniform distribution with larger sample sizes.

**Poisson: distribution P-values**

**Figure S3.1.** Distribution of p-values for the Poisson GLMs for each dispersion test.
10,000 simulations per simulation set (intercept x sample size).

**Figure S3.2.** Distribution of p-values for the binomial GLMs for each dispersion test. 10,000 simulations per simulation set (intercept x sample size).

## S4. Dispersion statistics for GLMs

129

130   The dispersion statistics of the tests for GLMs tended to be smaller than 1

131 (expected value) when there was no overdispersion simulated for very small sample

132 sizes for both binomial and Poisson distributions (Figure S4.1). The exception was the

133 nonparametric Pearson test that presented values larger than 1 for the very small

134 intercepts (-3 in both distributions, 3 in binomial only). When comparing dispersion

135 statistics for the simulated overdispersed data (Figures S4.2 and S4.3), we found that

136 both Pearson-based dispersion statistics presented similar values. In contrast, the

137 dispersion statistic of the simulation-based residual variance presented lower values for

138 small sample sizes. The differences in dispersion statistics between tests tended to

139 increase with the increase of simulated overdispersion, but in opposite directions for

140 binomial and Poisson GLMs (Figure S4.4 and S4.5). Moreover, we found out that the

141 dispersion statistics of the simulation-based residual variance test depend heavily on the

142 slope parameter of the simulated data (Figure S4.6).

143

**Figure S4.1.** Median of the dispersion statistics of the tests for A) Poisson and B)
binomial GLMs, simulated without overdispersion for different intercepts (panels) and
sample sizes (x-axis) for the three dispersion tests: parametric Pearson test,
nonparametric Pearson test, and simulation-based residual variance test. The dotted
horizontal line indicates the ratio of 1. Values below the line are considered
underdispersion, and above the line are overdispersion. For all simulations, the slope
was fixed at 1.

**Figure S4.2.** Dispersion statistics (median) for GLM Poisson. Notice the different y-axis scales across sample sizes.

**Binomial: dispersion statistics**
1000 sim; Ntrials=10

test — Sim-based residual variance — param. Pearson residuals — nonparam. Pearson residuals

154

155  **Figure S4.3.** Dispersion statistics (median) for GLM binomial.

156

**Figure S4.4**. The dispersion statistics of the simulation-based residual variance test are smaller than the parametric Pearson test statistics for all binomial models and for small sample sizes in Poisson models. The differences between the two dispersion statistics decrease with increasing sample size (coloured lines) and increase with simulated overdispersion in the data (x-axis). The relative differences (y-axis) were calculated by subtracting the simulation-based dispersion statistics from the parametric Pearson statistic, then dividing by the simulation-based statistic, and can be interpreted as the difference in the percentage of the simulation-based statistics. The results presented are based on 1,000 simulations with zero intercepts.



166

**Figure S4.5.** The dispersion statistics of the simulation-based residual variance test are smaller than nonparametric Pearson dispersion statistics for all binomial models and for small sample sizes in Poisson models. The differences between the two dispersion statistics decrease with increasing sample size (coloured lines) and increase with simulated overdispersion in the data (x-axis). The relative differences (y-axis) were calculated by subtracting the Parametric Bootstrapping statistics from the simulation-based dispersion statistics, then dividing by the simulation-based statistics, and can be interpreted as the difference in the percentage of the simulation-based statistics. The results presented are based on 1,000 simulations with zero intercepts.

15

## S5: Expanding simulation parameters for GLMs

176

177    Here, we investigated the possible influence of other parameters used to generate

178    the datasets for binomial and Poisson GLMs. In Figure S5.1, we investigated the power

179    and dispersion statistic for datasets simulated with different slopes (the default slope in

180    all other simulations was 1). In Figure S5.2, we investigated the effect of varying the

181    number of trials on the binomial GLMs in terms of power, type I error, and dispersion

182    statistics.



183

**Figure S5.1**. Power and dispersion statistics for simulations with different slopes (panel
columns) for binomial and Poisson GLMs. Number of simulations = 500; intercept = 0,
number of trials for the binomial = 10.

187

**Figure S5.2.** Power, dispersion statistics, and type I error of dispersion tests for binomial data simulations with different numbers of trials (panel columns). The fixed parameters are: intercept = 0, sample size = 500, slope = 1. Results for 1000 simulations.

## S6. Power for the GLMs

*Power calibration*

To investigate if the lower power of the simulation-based residual variance test is a consequence of the very conservative type I error rates, we calibrated the power using the p-value at the 5% quantile of the empirical distribution of p-values where the null hypothesis was true for each set of simulations (Figures S3.1 and S3.2). This method should provide an estimate of differences in power, controlling for type I error rate (Luke et al. 2017). Figures S6.1 and S6.2 show the power (calibrated and uncalibrated) of the dispersion tests for each simulation set (intercept, sample size and overdispersion) for Poisson and binomial GLMs, respectively.

**Figure S6.1.** Power for GLM Poisson.

**Binomial power**

204

**Figure S6.2.** Power for GLM binomial.

206 **S7. Additional GLMM results**

*Type I error rate of the alternative dispersion tests*

208        In Figures S7.1 and S7.2, we present the type I error rates for the four alternative

209   dispersion tests for the Poisson and binomial GLMMs, respectively, using simulated

210   sets of parameters: number of observations, number of groups, and intercepts.



211

**Figure S7.1.** Type I error rate for the three alternative dispersion tests for the Poisson
GLMMs. 1000 simulations for each parameter set. To improve visualisation of the
different intercept lines, the x-axis values were slightly displaced to align with the
sample size values.

216

**Figure S7.2.** Type I error rate for the three alternative dispersion tests for binomial GLMMs. 1000 simulations for each parameter set. To improve visualising the different intercept lines, the values in the x-axis were slightly displaced around the sample size values.

*Power of the alternative dispersion tests*

In Figures S7.3 and S7.4, we show the Power for the three alternative dispersion tests for the Poisson and binomial GLMMs, respectively, for the simulated sets of parameters: number of observations, number of groups, and intercepts.

**Figure S7.3.** Power of the three alternative dispersion tests for the Poisson GLMMs, with different sample sizes (rows), intercepts (columns), and number of groups for the random intercept (line types). The missing lines for the first panel (intercept = -3 and sample size = 50 are due to simulation errors for some tests. For each parameter set, we ran 1000 simulations.
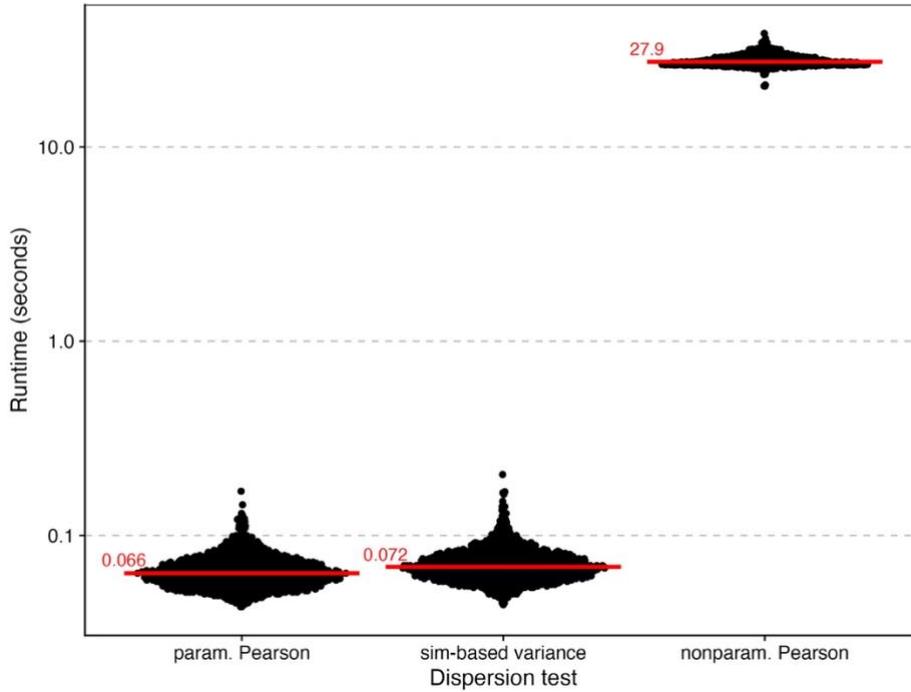
**Figure S7.4.** Power of the three alternative dispersion tests for binomial GLMMs, with different numbers of observations (rows), intercepts (columns), and number of groups for the random intercept (line types). 1000 simulations for each parameter set.

*Dispersion statistics of the alternative dispersion tests*

In Figures S7.5 and S7.6, we show the dispersion statistics for the three alternative dispersion tests for the Poisson and binomial GLMMs, respectively, for the simulated sets of parameters: number of observations, number of groups, and intercepts.

24

**Poisson: dispersion statistics**

239

**Figure S7.5.** Dispersion statistics of the three alternative dispersion tests for the Poisson GLMMs, with different numbers of observations (rows), intercepts (columns) and number of groups for the random intercept (line types). The missing lines for the first panel (intercept = -3 and sample size = 50 are due to simulation errors for some tests. For each parameter set, we ran 1000 simulations.

**Figure S7.6.** Dispersion statistics of the three alternative dispersion tests for binomial
GLMMs, with different numbers of observations (rows), intercepts (columns), and
number of groups for the random intercept (line types). 1000 simulations for each
parameter set.

*Computational runtime for tests with GLMMs*

We computed the run time for the three tests used for GLMMs: the parametric
Pearson test, the nonparametric Pearson test, and the simulation-based residual variance
test with conditional simulations (Figure S7.7). We used 1,000 simulations of the
Poisson GLMM as an example, with an overdispersion parameter of 0.4, an intercept of
0, a sample size of 1,000, and 100 groups. There was almost no difference in

26

256      computational time between the parametric Pearson test (median at 0.066 seconds) and

257      the simulation-based residual variance test (median at 0.072 seconds). As expected, the

258      nonparametric Pearson residuals presented the largest runtime, with a median of 27.9

259      seconds.



260

261 **Figures S7.7.** Runtime (in seconds) for each dispersion test for a Poisson GLMM
262 simulated 1000 times with the following parameters: overdispersion parameter of 0.4,
263 an intercept of 0, a sample size of 1,000, and a number of groups of 100. Note the y-axis
264 at the log 10 scale.

265

## S8: Alternative simulation-based residual variance test

Another possibility for improving dispersion tests for GLMMs is to develop a simulation-based approach that shows better type I, power, and a dispersion statistic that could be interpreted similarly to the Pearson dispersion. To explore future possibilities, we briefly considered an alternative simulation-based test that attempts to approximate the Pearson residuals by dividing the observed raw residuals (observed – fitted values) by the variance of the simulated values for each observation (Equations S8.1 and S8.2). We evaluated and compared this test for Poisson and binomial GLMs and GLMMs (conditional simulations only), as we did for the other tests.

$$Approx. Pearson\ observed\ residuals:\ r_i = \frac{(y_i - \hat{\mu})}{var(y_{is})} \qquad \text{(Equation S8.1)}$$

$$Approx. Pearson\ simulated\ residuals:\ r_{is} = \frac{(y_{is} - \hat{\mu})}{var(y_{is})} \qquad \text{(Equation S8.2)}$$

One obstacle with calculating the denominator of the approximate Pearson residuals for each observation is that the variance depends on the number of simulations and the model parameters, such as the intercept or the number of trials in the binomial GLM/GLMMs. If there are too few simulations or the intercept is very small, the chance of resulting in zero variance (all simulated values are the same) is higher for data points with small variance. To overcome this, we first evaluated the minimum number of simulations for different intercepts and sample sizes, in which all observations have estimated variances that are different from zero. For all combinations of parameters, we found that 1,000 simulations were sufficient to ensure that all variances in the simulated observations were positive (Figures S8.1 and S8.2). However, 250 simulations (the default parameter of the DHARMa package) also presented reasonable results, with the only exception being the Poisson GLMs with 30 out of 1,000 simulations (sample size

289    of 100 and intercept of -1.5) with a very low percentage of zero variances in the

290    simulated observations (mean of 0.01, maximum of 0.06). We are aware that the

291    number of zero variances in the simulations depends heavily on the simulation set, e.g.,

292    the number of trials for the binomial GLM. To develop an effective dispersion test, one

293    should consider alternatives to address this issue. For the subsequent analyses, we

294    excluded the simulations with zero variance in any simulated observation to compare

295    the alternative dispersion test with the simulation-based residuals test and the Pearson

296    Chi-squared dispersion test.



297

<cue>298</cue> **Figure S8.1.** Poisson GLM: Proportion of observations with simulated zero variance in
<cue>299</cue> the dataset for different combinations of intercept (columns), number of simulations
<cue>300</cue> (rows) and sample sizes (colours).



Binomial: Prop obs with zero SD simulated data

<cue>301</cue>
<cue>302</cue> **Figure S8.2.** Binomial GLM: Proportion of observations with simulated zero variance
<cue>303</cue> in the data set for different combinations of intercept (columns), number of simulations
<cue>304</cue> (rows) and sample sizes (colours). The number of trials of the binomial was set to 10 in
<cue>305</cue> all simulations.

<cue>306</cue>    First, we compared the approximate Pearson residuals for GLMs with the

<cue>307</cue> Pearson residuals by regressing the difference between them as the response variable

<cue>308</cue> and the Pearson residuals as the predictor for the Poisson GLMs (Figure S8.3). The

<cue>309</cue> intercepts for all simulation sets were nearly zero. The slope of the regression was

<cue>310</cue> positive and very small for the larger number of simulations and intercepts. It means

<cue>30</cue>

311 that the approximate Pearson tends to be slightly larger than the Pearson for larger

312 residuals.



313

**Figure S8.3.** Mean slope (A) and intercept (B) of the regression of the difference
between the Approximate Pearson residuals and Pearson residuals as response variable
and the Pearson residuals as predictor for the Poisson GLMs.

317        Type I error rates for the alternative simulation-based test, based on the

318 approximate Pearson residuals for GLMs, were similar to those for the simulation-based

319 residual variance test for the Poisson model. They tended to be conservative for small

320 intercepts (Figure S8.4). However, for the binomial model, type I error rates were more

321 similar to the parametric Pearson residuals test, with values closer to 0.05 (Figure S8.4).

322



323

**Figure S8.4**. Type I error rates for GLMs comparing the parametric Pearson residuals tests, the simulation-based residual variance test and the simulation-based approximate Pearson test.

The dispersion statistics for the alternative simulation-based residual variance test didn't change depending on the number of simulations and were very similar to the parametric Pearson dispersion statistics for both GLMs (Figure S8.5). Power was very similar among the tests for the Poisson GLM (Figure S8.6). For binomial GLMs, the power of the alternative simulation-based residual test was high and similar to the parametric Pearson residuals test.

333

334 **Figure S8.5.** Dispersion statistics GLMs. Simulation set with intercept = 0.



335

336 **Figure S8.6.** Power GLMs. Simulation set with intercept = 0.

337       For the GLMM simulations, we fixed the number of groups at 100 and the

338 number of simulations at 250 to compare with the cases where the Pearson Chi-squared
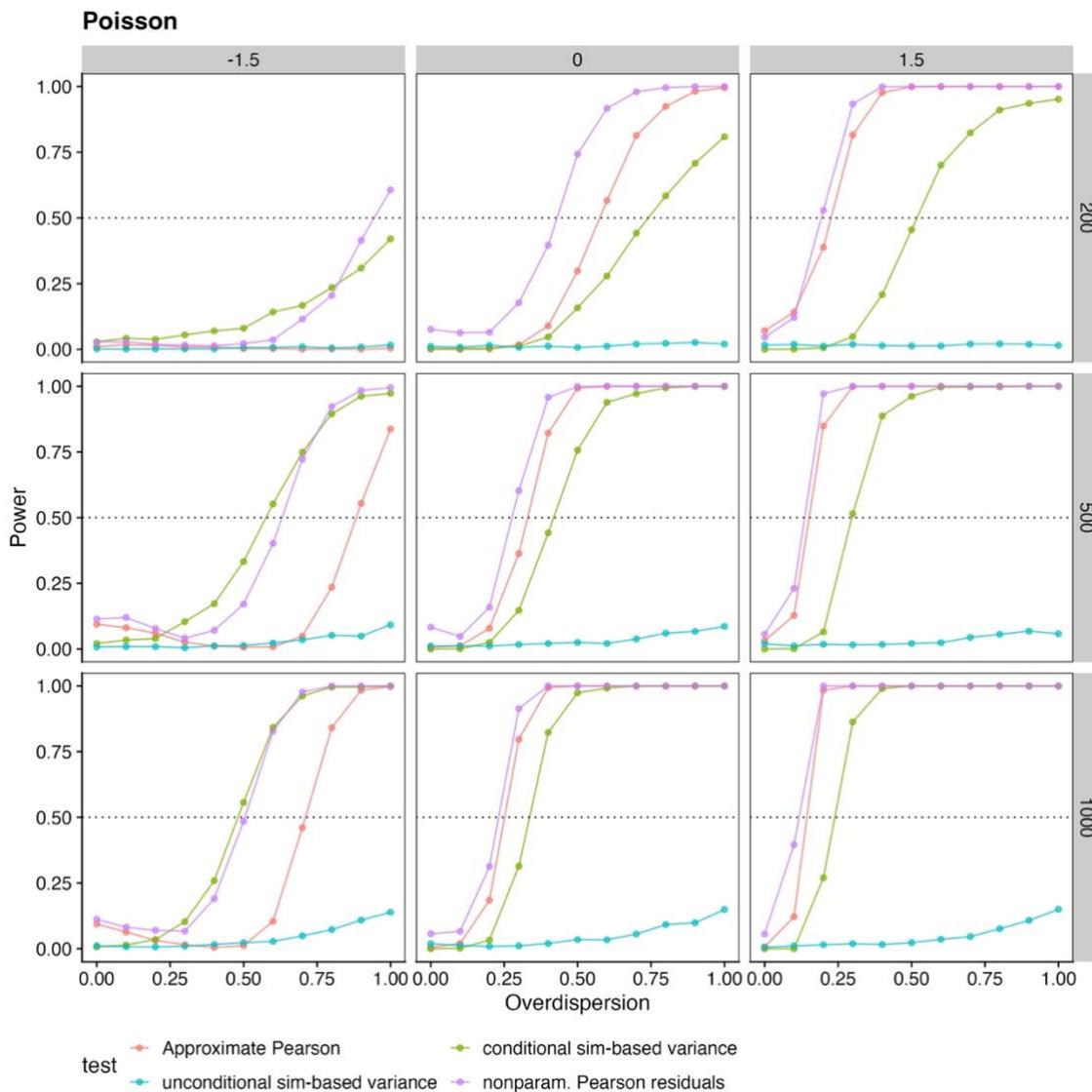
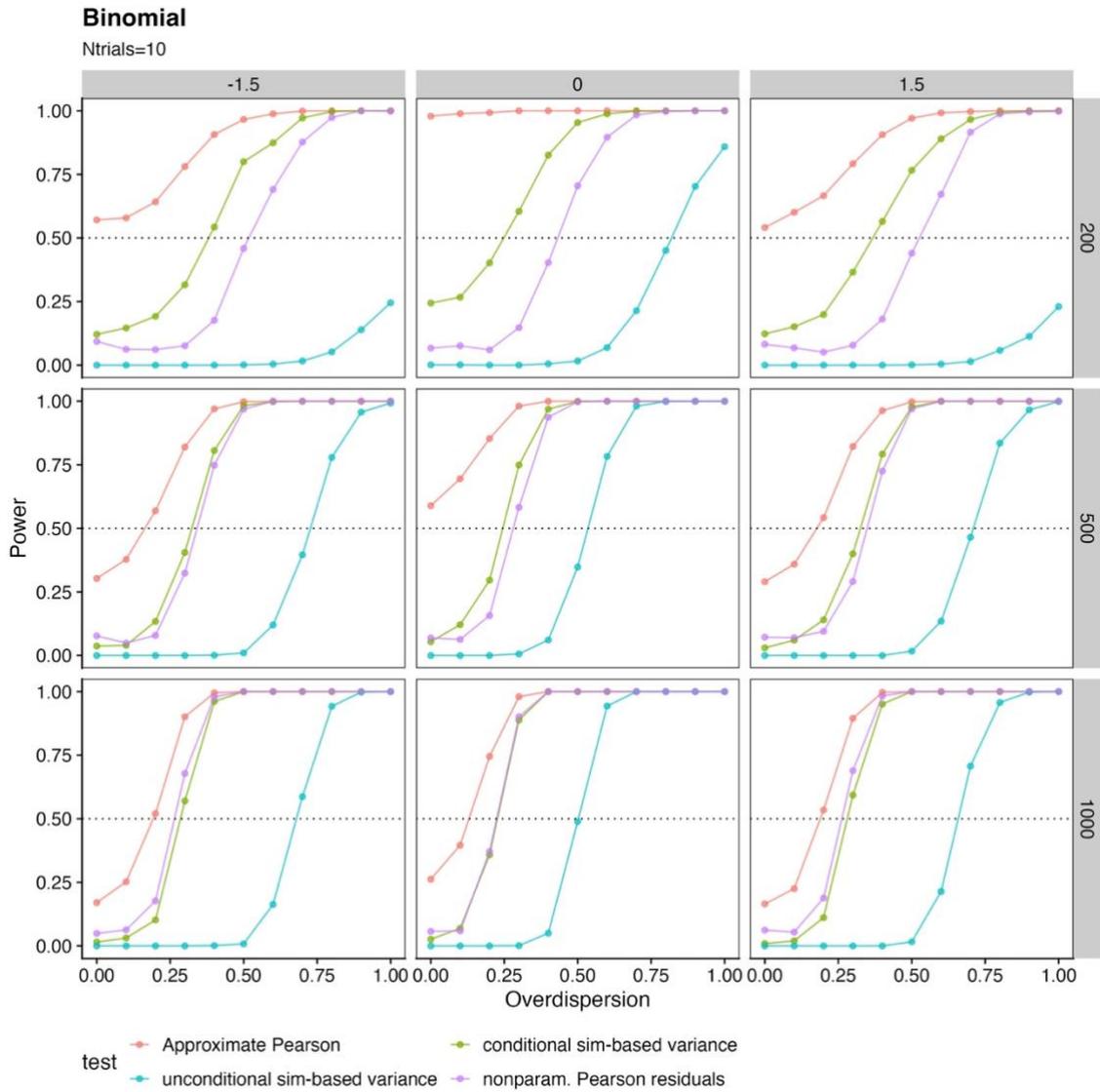339 test fails. We compared sample sizes of 200, 500, and 1000 observations and intercepts

340 of -1.5, 0, and 1.5. We excluded simulations with zero variance in the simulated

341 observations (specifically, for Poisson GLMMs, which accounted for less than 0.1% of

342 the simulations). For GLMMs, we used only the conditional simulations, which have

343 been proven to yield better dispersion test results.



344

345 **Fig S8.7.** Power for Poisson GLMMs for the alternative simulation-based test using an
346 approximation for Pearson residuals compared with the other tests assessed in the study.
347 1000 simulations for each parameter set: intercept (panel columns) and sample size
348 (panel rows). The fixed parameters are slope = 1, number of groups = 100, and random
349 effects variance = 1.

**Fig S8.8.** Power for binomial GLMMs for the alternative simulation-based test using an approximation for Pearson residuals compared with the other tests assessed in the study. 1000 simulations for each parameter set: intercept (panel columns) and sample size (panel rows). The fixed parameters are slope = 1, number of groups = 100, random effects variance = 1, number of trials = 10.

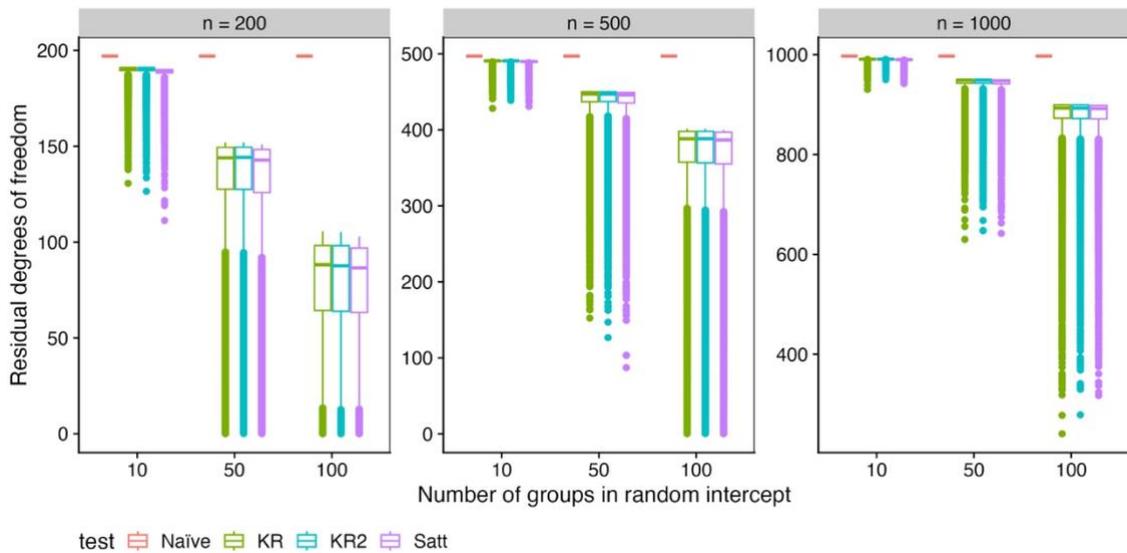## S9. Parametric Pearson test with approximated residual degrees of freedom for GLMMs

Degrees of freedom (*df*) are not always known for GLMMs with complex hierarchical structures and limit the use of the parametric Pearson test because it depends on it for evaluating overdispersion with the Chi-squared distribution. Moreover, our results show that using the naïve *df* is problematic for testing dispersion when you have a large number of groups in the random intercept. The two most suggested methods to approximate *df* of mixed-effect models, the Satterthwaite (1946) and the Kenward-Roger (Kenward & Roger 2009), were developed for LMMs to account for the effect of the covariance structure on *df* and standard errors. Stroup et al. (2013) suggested that the adjustment is also accurate for GLMMs. However, none of the most used R packages use any correction for the degrees of freedom for GLMMs. The few R packages that provide those approximations, e.g. *lmerTest* (Kuznetsova et al., 2017; Kuznetsova et al., 2020) that relies on *pbkrtest* (Halekoh & Højsgaard 2014), are only implemented for LMMs.

Recently, we found that the R package *glmmrBase* (Watson 2024) provides those approximation methods for GLMMs. Thus, we compared the parametric Pearson test with the three corrections for degrees of freedom available in the package for the Poisson GLMMs. The corrections are:

- The Kenward-Roger (KR) bias-corrected variance-covariance matrix for the fixed effect parameters and degrees of freedom from Kenward & Roger (1997).
- The improved correction of the Kenward-Roger (KR2) returns an improved correction given in Kenward & Roger (2009).
- The Satterthwaite correction (Sat) from Satterthwaite (1946).
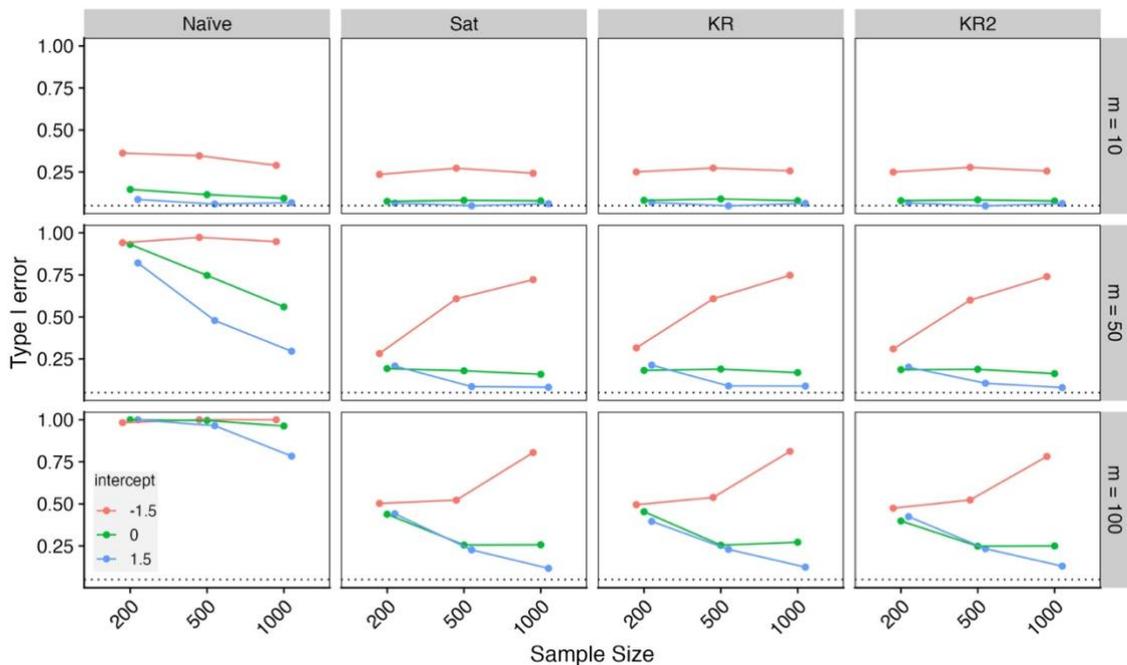
380        Our test results show that all three correction methods presented very similar

381    residual *df* for all simulation settings (Figure S9.1), which resulted also in very similar

382    test results (e.g., Figure S9.2 for type I error). Given the high similarity among tests for

383    the different residual *df* corrections, we show and discuss the results for the KR2 test in

384    comparison with the parametric Pearson "naïve" test and the alternative GLMM tests

385    (nonparametric Pearson and simulation-based residual variance test with conditional

386    simulations). In Figure S9.3, we observe that the correction for the residual *df* corrected

387    the dispersion statistics towards 1 for simulations without overdispersion, except for the

388    very small intercept (-1.5). This results in the two-sided dispersion test being less prone

389    to being significant, given the very low dispersion parameter (detecting underdispersion

390    instead of overdispersion).

391        Although the parametric Pearson tests with the approximated residual degrees of

392    freedom performed much better than those with the "naïve" residual *df*, they still

393    underperformed compared to the nonparametric version when having a large number of

394    groups in the random effects (Figure S9.4), especially for very small intercepts and
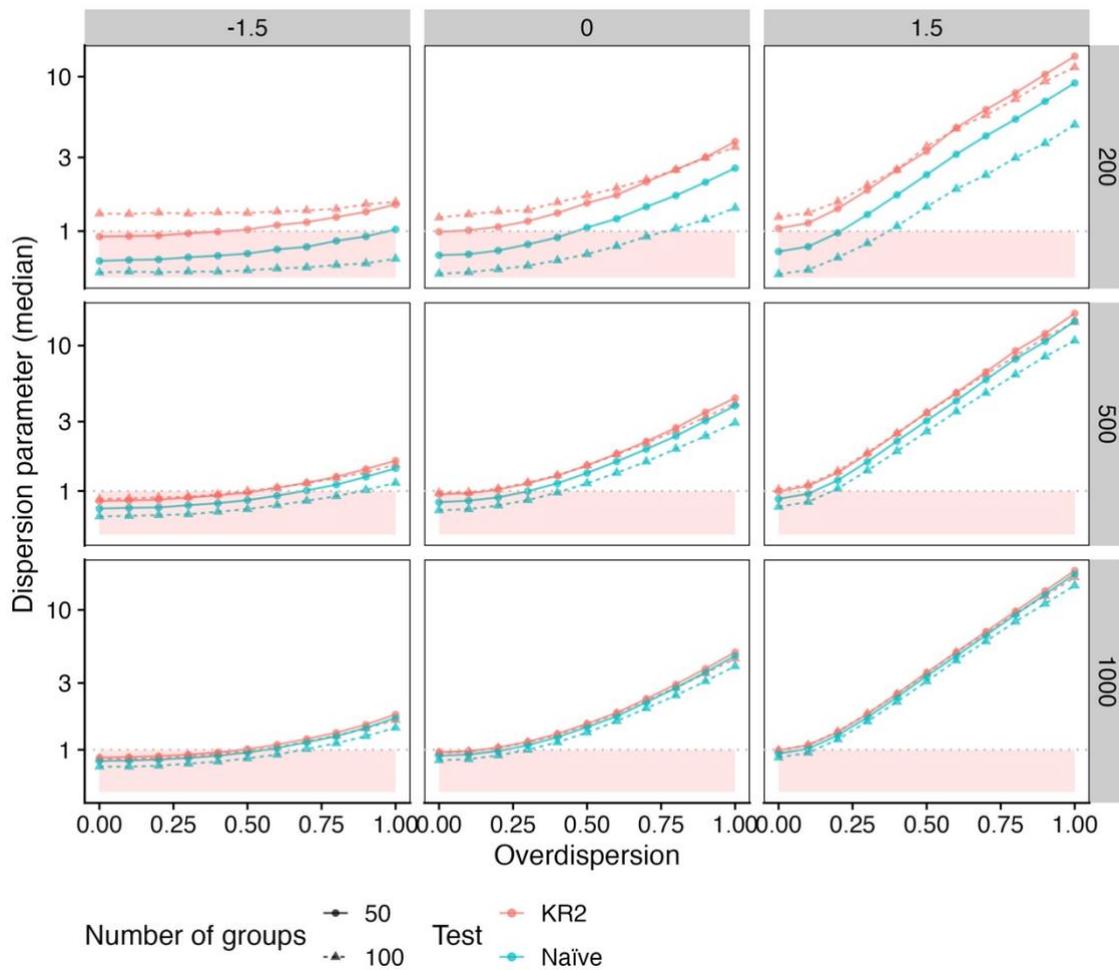
395    sample sizes.

396

**Figure S9.1.** Residual degrees of freedom for the different correction methods for
Poisson GLMMs with different numbers of groups in the random intercept (x-axis) and
sample sizes (panel columns). Please refer to the main text above to relate to each
applied correction. 1,000 simulations for each parameter setting, slope = 1, random
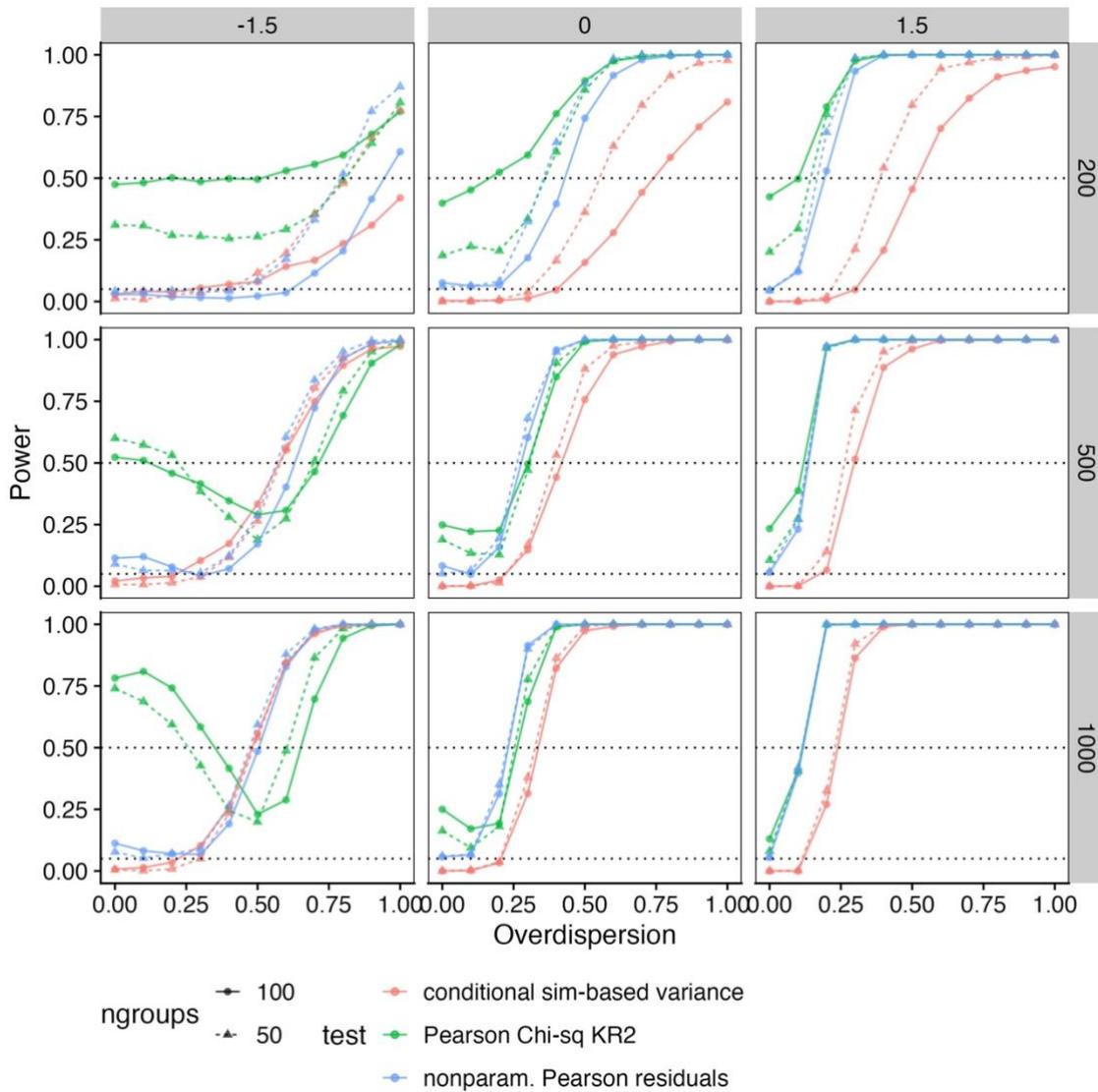intercept variance = 1.



402

**Figures S9.2.** Type I error for the parametric Pearson test for Poisson GLMMs
performed with different corrections for the residual degrees of freedom (panel
columns), number of groups in the random intercept (panel rows) and sample size (x-
axis). Data were simulated from a Poisson GLMM with different intercepts (colours).
Please refer to the main text above to relate to each applied correction. 1000 simulations
for each parameter setting, slope = 1, random intercept variance = 1.

**Figure S9.3**. Dispersion parameters for the parametric Pearson test for Poisson GLMMs performed with different corrections for the residual degrees of freedom (colours), number of groups in the random intercept (linetype and shape), sample size (panel rows), and intercept (panel columns). Please refer to the main text above to relate to each applied correction. To improve clarity, we omitted the other corrections because they are too similar to each other. 1000 simulations for each parameter setting, slope = 1, random intercept variance = 1.

417

**Figure S9.4**. Power of dispersion tests for Poisson GLMMs (colours) performed with
different numbers of groups in the random intercept (linetype and shape), sample size
(panel rows), and intercept (panel columns). Please refer to the main text above to relate
to the applied correction for residual degrees of freedom. To improve clarity, we omitted
other corrections for residual degrees of freedom because they are too similar to each
other. 1000 simulations for each parameter setting, slope = 1, random intercept variance
= 1.

# References

Hartig, F. (2024). *DHARMa: Residual Diagnostics for Hierarchical (Multi-Level / Mixed) Regression Models* (Version 0.4.7) [Computer software]. https://CRAN.R-project.org/package=DHARMa

Jahn, N. (2023). *europepmc: R interface to the europe PubMed central restful web service* (Version 0.4.3) [Computer software]. https://CRAN.R-project.org/package=europepmc

Laumer, I. B., Kansal, S., Van Cauwenberghe, A., Rahmaeti, T., Setia, T. M., Mundry, R., Haun, D., & Schuppli, C. (2025). Wild and zoo-housed orangutans differ in how they explore objects. *Scientific Reports*, *15*(1), 14853. https://doi.org/10.1038/s41598-025-97926-z

Lüdecke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021). performance: An R package for assessment, comparison and testing of statistical models. *Journal of Open Source Software*, *6*(60), 3139. https://doi.org/10.21105/joss.03139