

2 **Dispersion tests in generalized linear mixed-effects models: A**
3 **methods comparison and practical guide for ecologists**

4 Melina de Souza Leite^{1*}, Daniel Rettelbach^{1,2} & Florian Hartig¹

5 1. Theoretical Ecology, University of Regensburg, Germany

6 2. coTrial Associates, Department of Surgery, University Hospital Regensburg,
7 Germany (current address)

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

9 **Open research statement**

10 The data, code for the simulations, analyses, and figures presented in this manuscript
11 are available on Zenodo (<https://doi.org/10.5281/zenodo.17611060>) (Leite, 2026).

12 **Abstract**

13 Underdispersion and overdispersion are common issues when analysing ecological data
14 with generalized linear (mixed) models (GLMs/GLMMs). Overdispersion, the
15 phenomenon where observations spread wider than expected by the fitted model,
16 usually leads to anti-conservative p-values and, thus, to inflated type I error. In contrast,
17 underdispersion, a narrower spread of the data than expected, causes overly
18 conservative p-values and, therefore, reduced power. A range of tests has been proposed
19 to detect such dispersion problems, but there are few comparative studies of their
20 performance across models and analysis settings, and, most importantly, sparse
21 recommendations for ecologists on how to check for dispersion issues. Our goal was to
22 identify a general dispersion test for GLMs/GLMMs applicable to standard distributions
23 and random-effects structures commonly used in ecological data analysis. Following an

24 initial review of available tests, we selected two classes of dispersion tests: (1)
25 parametric and nonparametric tests based on Pearson residuals and (2) simulation-based
26 tests that compare the expected and observed residual variance. Comparing their
27 performance by type I error, power, and dispersion estimate, across a range of Poisson
28 and binomial GLMs/GLMMs, we found that the nonparametric Pearson residuals test
29 performed best across all metrics, particularly for data with low incidence or count rates
30 and/or small samples; however, at the cost of high computational expense. The
31 parametric Pearson residuals test, recommended in many books and guidelines, was fast
32 and effective for GLMs, but biased towards underdispersion in GLMMs due to the
33 naïve computation of the random-effect degrees of freedom. The simulation-based
34 residual variance test was slightly less powerful, but showed overall good calibration.
35 The latter offers a compromise between the strengths and weaknesses of the two
36 Pearson-based tests. We conclude that for GLMs, the parametric Pearson residuals test
37 offers the best balance of speed and accuracy. For GLMMs, we recommend either the
38 computationally demanding nonparametric Pearson residuals test or the faster, although
39 somewhat less powerful, simulation-based residual variance test. We also analyze two
40 case studies in ecology that differ in complexity and include recommendations for
41 ecological data analysis to address dispersion issues, using the most commonly used R
42 packages, avoiding pitfalls, and improving model fit and the interpretation of ecological
43 datasets.

44 **Keywords:** overdispersion/underdispersion, multilevel/hierarchical models, hypothesis
45 test, Pearson residuals, type I error, power, dispersion parameter

46 **Introduction**

47 Generalized linear models (GLMs) and generalized linear mixed models (GLMMs) are
48 the most commonly used tools for the statistical analysis of ecological data (Bolker et
49 al., 2009; Lai et al., 2019; Touchon & McCoy, 2016). By incorporating mixed and
50 random effect structures with a wide array of distributional assumptions (e.g., binomial,
51 Poisson), GLMMs allow researchers to model nonnormal response variables (e.g.,
52 counts, proportions, or presence-absence) while properly accounting for variation
53 clustered in sampling units, sites, or study years (Bolker et al., 2009; McMahon & Diez,
54 2007). However, as for all parametric statistics, these models rely on the fact that
55 residuals scatter around the regression mean with the specified distribution, and their
56 inferential results can be seriously biased if these distributional assumptions are
57 violated.

58 A particularly common and dreaded violation of distributional assumptions in
59 GLMs/GLMMs is overdispersion. Overdispersion refers to greater variation in the
60 observed data (and particularly the model residuals) than the fitted model assumes (H.
61 Campbell, 2021; McCullagh & Nelder, 1989). Strong overdispersion usually appears in
62 distributions that assume a fixed mean-variance relationship, such as the Poisson model
63 for count data (Harrison, 2014; J. M. Hilbe, 2014) or the binomial model for discrete
64 proportions (Dunn & Smyth, 2018; Harrison, 2015). For example, a Poisson process
65 assumes that we count randomly distributed points in space. However, when individuals
66 are subject to spatial/temporal clustering due to different ecological mechanisms (e.g.,
67 patchy resource distribution, social behaviour, dispersal limitation) and/or imperfect
68 detection (Rhodes, 2015), we typically find higher dispersion than expected from a
69 Poisson distribution (Box 1). Alternatively, overdispersion may also arise from

70 modeling misfit, for example, by failing to include important predictors and interactions
71 or by specifying the incorrect link function (J. M. Hilbe, 2011).

72 Overdispersion is a major concern in practical data analyses because it can have
73 substantial anti-conservative effects on p-values, confidence intervals, and all other
74 goodness-of-fit and precision metrics (Fig. 1, see also Rhodes, 2015). Anti-conservatism
75 means that p-values and confidence intervals are too small, leading to inflated false-
76 positive results (type I errors). In practice, we have encountered analyses where an
77 overdispersed model had very small and significant p-values (<0.001) that became
78 nonsignificant after switching to a GLM with more appropriate dispersion (see example
79 in Fig. 1).

80 The counterpart to overdispersion is underdispersion, where the variation in the
81 observed data (and, thus, model residuals) is lower than assumed by the fitted model.
82 Reasons for underdispersion can again be that the data-generating process (e.g., a
83 uniform distribution of individuals in space, Box 1) differs from what is assumed by the
84 model (Lynch et al., 2014). However, in practice, it is often the result of model
85 overfitting, i.e., having a overly complex model that overfits the data. Underdispersion
86 is somewhat less discussed in the ecological literature, both because it is less frequent,
87 but also because it leads to over-conservative model metrics (Fig. 1). This may seem
88 less problematic as it does not lead to reporting “wrong” effects, but underdispersion
89 reduces overall power and thus increases type II error. Therefore, accurate statistical
90 inference demands that we identify and adequately address both underdispersion and
91 overdispersion to minimise the risk of wrong inference.

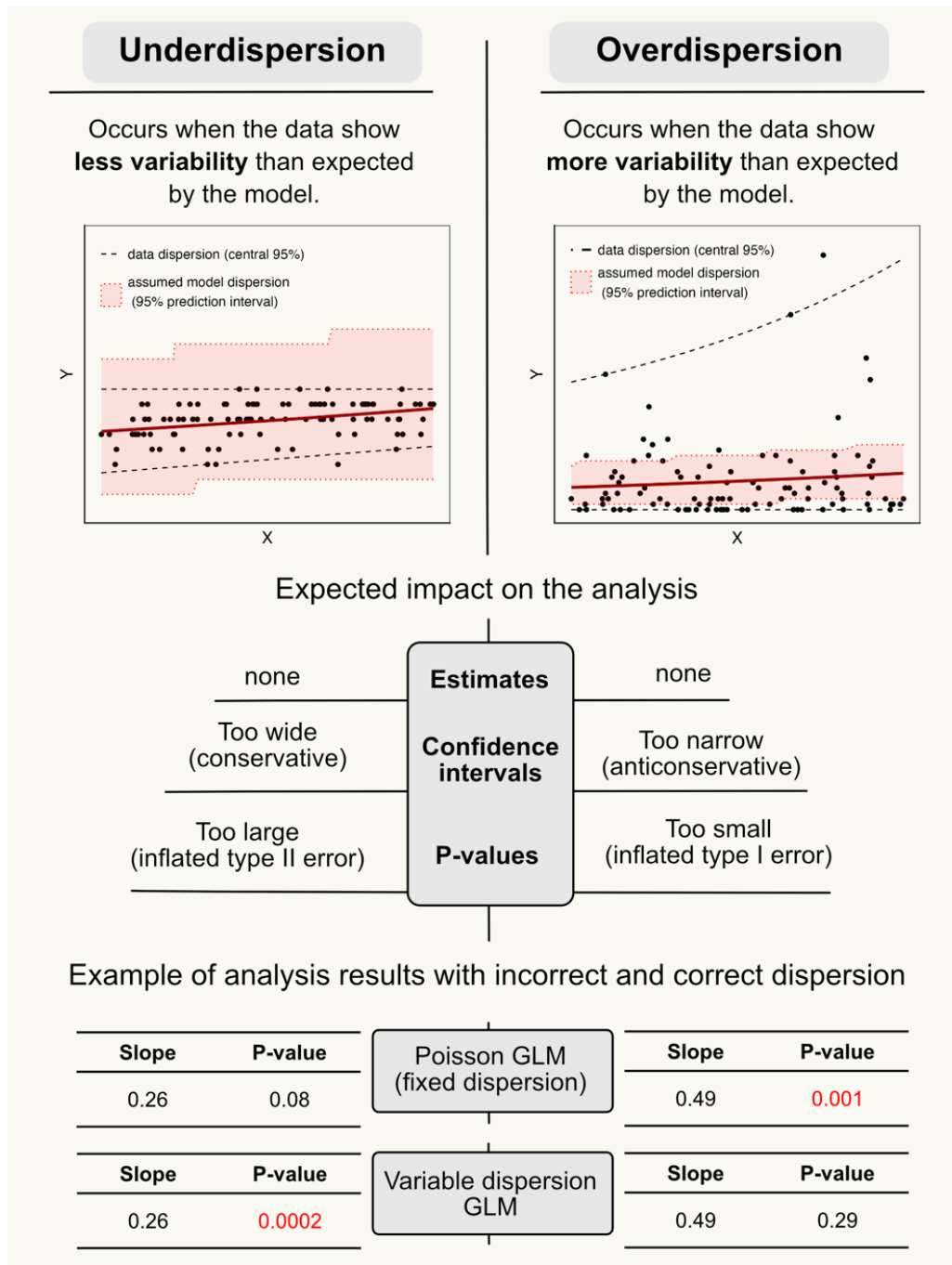
92 Given the central importance of dispersion to all statistical indicators,
93 statisticians have pondered how to detect and address dispersion problems since the
94 early days of modern statistics in the 19th century (see Quine & Seneta (1987) and

95 Xekalaki (2014) for a historical perspective). Since then, a large variety of approaches
96 have been proposed and discussed to deal with the “dispersion problem”, ranging from
97 (1) comparing models with or without free dispersion parameters through likelihood
98 ratio test, such as Poisson and negative binomial (e.g. Yang et al., 2007), (2) designing
99 specific hypothesis tests for the “extra” variation (e.g. Fisher, 1950), such as score tests
100 (Dean, 1992; Dean & Lawless, 1989; Lawless, 1987), (3) using goodness-of-fit tests,
101 such as tests on Pearson or Deviance residuals (Dunn & Smyth, 2018; McCullagh,
102 1985) (although the distinction between categories (2) and (3) can be blurry, see
103 (Collings & Margolin, 1985; Dean, 1992; Dean & Lawless, 1989) or (4) using
104 simulation-based nonparametric tests to compare observed and predicted variance of the
105 residuals (Hartig, 2026).

106 Somewhat confusing for the ecological data analyst, however, many of these
107 approaches have been designed and tested only in very specific scenarios (e.g., only for
108 a Poisson GLM), and there is a surprising lack of systematic evaluation of these tests
109 and strategies across a range of more complex GLMMs. Moreover, a quick review of
110 current methods in the R environment (R Core Team, 2024) revealed that existing
111 dispersion tests are scattered across multiple packages (Table 1), and most apply only to
112 a restricted set of models. In the ecological literature, although awareness of dispersion
113 problems has increased over the last 20 years (Box 2), there is still a clear lack of
114 guidance on how dispersion problems are assessed or tested (Box 2). All this makes it
115 challenging to decide which test to use in applied ecological data analysis.

116 The goals of this study are: (1) to review and order the diversity of dispersion
117 tests for GLMs and GLMMs, (2) to identify tests that can reliably work across a range
118 of models with diverse distributions and complex hierarchical structures, common
119 situations for ecological data analyses, and (3) to raise awareness and offer

120 recommendations to the ecological community by presenting two case studies and
121 suggesting tools. Following our literature review (next section), we identified two
122 groups of tests that appeared to be generally applicable: parametric and nonparametric
123 tests on Pearson residuals, as well as a new simulation-based nonparametric test that
124 directly compares observed and predicted variance of the raw residuals. We then used
125 simulated data to compare the performance of these tests in terms of type I error, power,
126 and the interpretability of the dispersion statistics. Based on the simulation results and
127 the case studies in ecology, we provide recommendations for the most suitable tests for
128 detecting over- or underdispersion, depending on model complexity and software
129 availability, i.e., currently available R packages and functions.



130

131 **Figure 1.** Definition, statistical consequences, and a practical analysis example of
 132 under-/overdispersion in GLMs/GLMMs. The top row shows examples of a data
 133 analysis using a Poisson GLM with simulated under- and overdispersed count data.
 134 Data points in black are contrasted to the Poisson model's 95% prediction interval (in
 135 red). Black dashed lines illustrate the data dispersion (central 95% quantiles). In the
 136 example, we present slope estimates and p-values for the GLM Poisson models fitted to
 137 the under- and overdispersed data above, and the results from more appropriate models
 138 with correct dispersion: a Conway-Maxwell-Poisson GLM for underdispersed data and
 139 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 reflect 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 modeling processes can lead to deviations from this distribution, resulting in over- or underdispersion. For example, spatial aggregation (clustering) of individuals due to patchy resource distribution, social behaviour, or dispersal limitation increases sampling variability and thus creates overdispersion compared to the Poisson (Fig. B1, see also Lindén & Mäntyniemi, 2011). In contrast, 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.

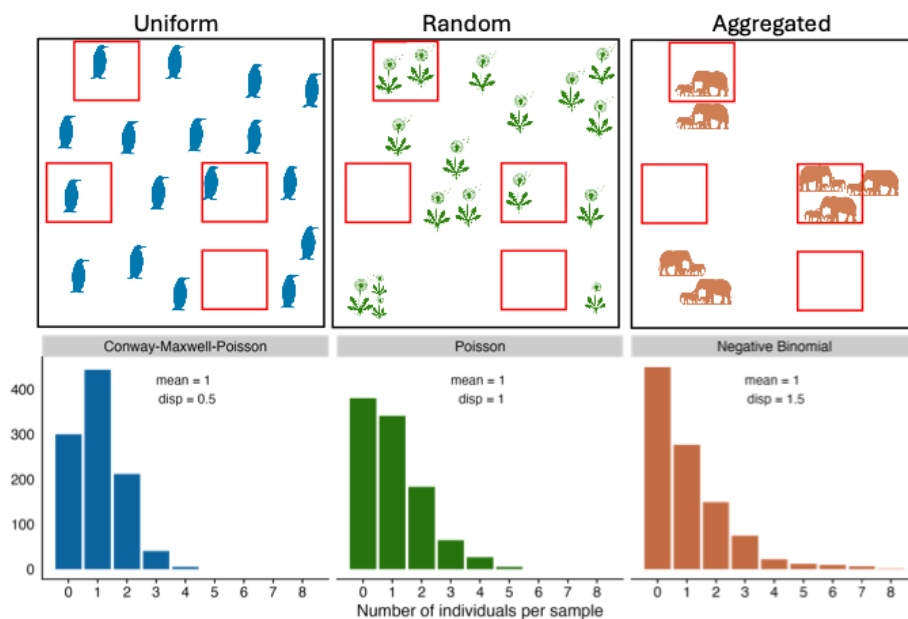


Figure B1. Examples of spatial distribution patterns of individuals under the same sampling design (quadrats in red) and the corresponding data-generating distributions for each pattern: Poisson (random pattern, green), negative binomial (aggregated pattern, light brown), and Conway-Maxwell-Poisson (uniform 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., N. A. Campbell & Reece, 2005): most penguin species are territorial, 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 constraints, such as ovule number, parental care or resources availability (M. E. 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 at 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 GLMs/GLMMs for such data has 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).

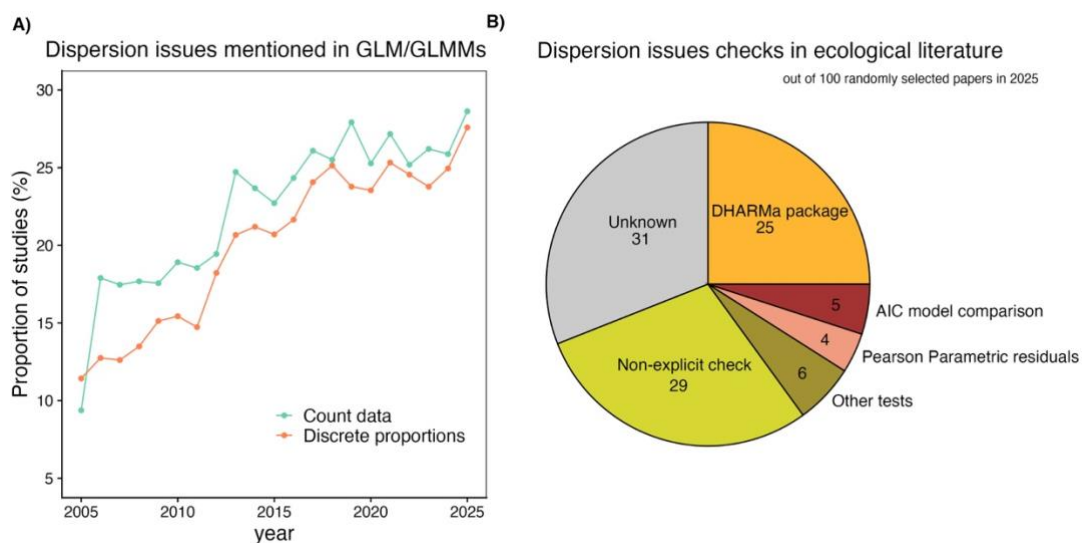


Figure B2. **A)** Annual trends for the proportion of ecological studies using GLMs/GLMMs for count and/or discrete proportion data **that mention dispersion terms** in the text. **B)** The type of checks and tools used for dispersion problems found in the papers that mentioned dispersion terms in 100 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 is 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: - Poisson and negative binomial or generalized Poisson - binomial and beta-binomial	<code>pscl::odTest()</code>	GLM Poisson -> negative binomial with <code>MASS::glm.nb()</code>	Jackman (2024)	
			<code>DCluster::test.nb.pois()</code>	GLM Poisson -> negative binomial with <code>MASS::glm.nb()</code>	Lopez-Quirez (2005)	
			<code>anova(..., test="LRT")</code>	Many GLMs/GLMMs*	R Core Team (2024)	
			<code>lmtest::lrtest()</code>	GLMs	Zeileis & 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.	<code>DCluster::DeanB()</code>	GLM Poisson.	Lopez-Quirez (2005)	
			<code>DCluster::DeanB2()</code>	Score tests based on Dean (1992)		
	Regression-based test for overdispersion from Cameron & Trivedi (1990)	Distribution specific (Poisson-based only).	<code>Rfast2::overdispreg.test()</code>	GLM Poisson (own model implementation)	Papadakis et al. (2025)	
			<code>overdisp::overdisp()</code>	GLM Poisson (own model implementation)	Cameron & Trivedi (2023)	
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. For complex models, difficult to define parametric null distribution (unclear residual degrees of freedom).	<code>msme::P__disp()</code>	GLMs	Hilbe & Robinson (2025)	
			<code>aods3::gof()</code>	GLMs	Lesnoff et al. (2024)	
			<code>DHARMA::testDispersion(..., type="Pearson")</code>	GLMs/GLMMs (naïve residual <i>df</i>)	Hartig (2026)	
			<code>performance::check_overdispersion()</code>	GLMs/GLMMs (naïve residual <i>df</i>)	Lüdecke (2021)	
			<code>RVAideMemoire::overdisp.glmmer()</code>	GLMMs (from <code>lme4</code> package, naïve residual <i>df</i> , calculates only dispersion statistic, no test)	Herve (2025)	
			<code>DHARMA::testDispersion(..., refit=T, type="Pearson")</code>	GLMs/GLMMs	Hartig (2026)	
			Deviance residuals: assumes the residual deviance are chi-squared distributed.	<code>aods3::gof()</code>	GLMs	Lesnoff et al. (2024)
			Dispersion metric: the square root of the penalized residual sum of squares divided by the number of observations.	<code>blmeco::dispersion_glmmer()</code>	GLMMs from <code>lmer::glmer()</code> (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. Fast nonparametrics but possibly less exact than working on the standardized residual dispersion.	<code>DHARMA::testDispersion(..., type="DHARMA")</code>	GLMs/GLMMs	Hartig (2026)	

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

144 **A short review of existing approaches to dispersion tests**

145 After reviewing the available literature, we divided the proposed strategies for
146 addressing dispersion problems into four classes (Table 1). Here, we discuss these broad
147 strategies in more detail and explain why we focused on two of these classes as the most
148 suitable competitors for a general dispersion test for GLMs and GLMMs. We note that, in
149 addition to the four approaches mentioned here, dispersion problems may also show up in
150 general goodness-of-fit tests (e.g., Feng et al., 2020). However, as they are not specifically
151 designed to react to dispersion, we did not consider them further.

152 *Likelihood ratio tests*

153 A first general strategy for detecting dispersion problems is to compare a model
154 with fixed dispersion to its nearest “relative” with variable dispersion using a likelihood
155 ratio test (LRT) or another model selection technique, such as AIC (Yang et al., 2007). For
156 count data, this could involve comparing a Poisson GLM to a negative binomial or
157 generalized Poisson GLM (J. M. Hilbe, 2014), or comparing a binomial GLM to a beta-
158 binomial GLM (Dunn & Smyth, 2018). While relatively easy to implement, the downside
159 of this approach, apart from the higher computational cost of fitting two models, is that it
160 doesn’t provide any direct diagnostics of over- or underdispersion. The alternative model,
161 however, might also fit better or worse for reasons other than a dispersion problem.
162 Moreover, using LRTs to detect dispersion problems has also been discouraged, as they
163 may yield unreliable results (Dean, 1992) and tend to underestimate the evidence against

164 the base model (Lawless, 1987). Therefore, we do not find this approach suitable as a
165 general dispersion test and do not consider it further.

166 *Score tests*

167 A second traditional option for assessing overdispersion is the score test (Dean,
168 1992; Dean & Lawless, 1989; Lawless, 1987). Score tests, also known as Lagrange
169 Multiplier (LM) tests, evaluate the gradient of the log-likelihood (called the score or LM
170 statistic) of a restricted parameter estimator (e.g., an overdispersion estimator constrained to
171 zero). Under the null hypothesis that the overdispersion is indeed zero, the score will have
172 an asymptotic chi-squared distribution (Rao, 1948). In performance comparisons, score
173 tests have been found to have good power (Ohara Hines, 1997), but their disadvantage is
174 that they are usually model specific (in the sense that different tests are needed for Poisson
175 or binomial GLMs); their implementation can be computationally demanding; and, as they
176 require access to the score, they must usually be implemented with the model and cannot be
177 calculated on top of a fitted model object. Perhaps because of these issues, we were unable
178 to find any R function that computes score tests beyond the Poisson GLM (Table 1),
179 although score tests have been developed for other models, such as the binomial GLM
180 (Dean, 1992).

181 An equivalent test related to the score test under certain conditions is the regression-
182 based overdispersion test proposed by Cameron & Trivedi (1990). Under a Poisson model,
183 the squared deviation of the observations from their fitted mean, after subtracting the
184 observation itself and scaling by the fitted mean, has expectation zero. In contrast, under
185 the negative binomial, it increases systematically with the mean. This motivates an

186 auxiliary regression of the transformed variable against the fitted mean, with a significant
187 slope indicating extra-Poisson variation. The main advantage of this test against other score
188 tests is its ease of implementation: it can be carried out after fitting a standard Poisson
189 GLM. However, similar to an LRT, the linear regression imposes a particular form of
190 overdispersion as an alternative hypothesis, and therefore, seems less general than the test
191 based on Pearson residuals described below.

192 We discarded score tests in general, and the Cameron & Trivedi (1990) test in
193 particular, from our further analysis, as it seems impractical to implement them across a
194 wide range of existing GLMM software.

195 *Tests based on residual dispersion*

196 A third class of testing approaches, arguably the most intuitive, directly calculates a
197 test statistic or goodness-of-fit metric on standardized model residuals. The most widely
198 used test of this kind is based on the sum of the model's Pearson residuals. As Pearson
199 residuals divide the raw residuals by the expected residual standard deviation, a correctly
200 specified model is expected to have a Pearson residual of around 1 for each observation. A
201 dispersion statistic is then defined as the sum of squared Pearson residuals divided by the
202 residual degrees of freedom. Models with a so-defined dispersion statistic > 1 are
203 considered overdispersed, while dispersion statistics < 1 are underdispersed. Sometimes,
204 this metric is modified by replacing the sum of squared Pearson residuals with the model
205 deviance, which is typically more readily available. However, as Venables & Ripley (2002)
206 discuss, this metric should be avoided, as it often deviates from 1, even for correctly
207 specified GLMs.

208 Defining dispersion via the Pearson statistic has the added advantage that for a
209 GLM, the expected distribution under the null hypothesis of a correctly specified model
210 asymptotically follows a chi-squared distribution (McCullagh, 1985). This allows a
211 straightforward construction of a hypothesis test, where we compare the Pearson statistic to
212 the chi-squared distribution with the respective residual degrees of freedom (*df*). This test is
213 referred to with different terminologies, such as the Pearson chi-squared dispersion test, the
214 Pearson residuals-based test for overdispersion, or simply the Pearson dispersion test.
215 Hereafter, we refer to this test as the **parametric Pearson residuals test** to differentiate it
216 from the nonparametric version discussed below.

217 An alternative approach to constructing a dispersion test based on the Pearson
218 dispersion statistic involves generating a null distribution through parametric bootstrapping.
219 A parametric bootstrap means that new data are simulated from the fitted model, and then
220 the statistic of interest (in this case: the Pearson statistic of a fitted model) is calculated
221 based on these data. The parametric bootstrap has been previously used for hypothesis tests
222 in mixed-effects models where parametric null distributions were difficult to obtain (e.g.,
223 Barr et al., 2013; Luke, 2017), and thus it seems a logical alternative for more complicated
224 models where the chi-squared distribution of the Pearson dispersion statistic cannot be
225 taken for granted (see methods for GLMMs below). Nevertheless, implementing parametric
226 bootstrapping in complex models can be less efficient for at least two reasons: it is time-
227 consuming and error-prone during model refits (Luke, 2017; Moral et al., 2017). A
228 dispersion test based on this principle was implemented in R by Hartig (2026). Hereafter,
229 we will refer to this test as the **nonparametric Pearson residuals test**.

230 *Tests based on simulated residual variances*

231 Simulation approaches can also be useful for generating null distributions for
232 alternative dispersion metrics. A final class of dispersion test approaches, which, to our
233 knowledge, was introduced in the DHARMA R package (Hartig, 2026) but has not been
234 discussed in the literature so far, involves defining a test statistic based on the dispersion of
235 the raw residuals. The test compares the observed raw residual variance (differences
236 between the observed data and the model predictions) with the simulated raw residual
237 variances (differences between the simulated data and the model predictions). Both
238 variances are scaled to the variance of all simulated observations to account for differences
239 in the number of simulations across the fitted model. For GLMMs, data simulations can be
240 generated conditionally or unconditionally on the fitted random effects. The dispersion
241 statistic is then defined as the ratio of the observed residual variance to the mean of the
242 simulated residual variances. Similar to the Pearson statistic, a ratio > 1 indicates
243 overdispersion, a ratio < 1 indicates underdispersion, and a significance test is constructed
244 based on the distribution of simulated residual variances.

245 From a theoretical perspective, this approach seems less elegant than the use of
246 Pearson residuals, because the latter, by “standardising” the residual dispersion relative to
247 the expected dispersion, allows each data point to contribute similarly to the dispersion
248 statistic. In contrast, the test on the unstandardised residuals will be more influenced by
249 large data points. However, the primary advantage of this approach is computational, as it
250 enables a nonparametric estimate of the test statistic without requiring a re-fit of the model
251 (in contrast to the nonparametric Pearson residuals test). Hereafter, we will refer to this test
252 as the **simulation-based residual variance test** to differentiate it from tests based on
253 Pearson residuals.

254 **Methods**

255 *Selected models and setup of the performance comparisons*

256 After reviewing the available approaches, we identified three tests as potential candidates
257 for a generally applicable dispersion test that could be implemented across a wide range of
258 GLMs and GLMMs:

- 259 (1) The parametric Pearson residuals test
- 260 (2) The nonparametric Pearson residuals test
- 261 (3) The simulation-based residual variance test

262 To compare the performance of these three tests, we simulated datasets based on the
263 two main distributions that often exhibit over- or underdispersion: the Poisson and the
264 binomial (N/K) discrete proportions. We varied the sample size (from 10 to 10,000) and the
265 intercept (from -3 to 3, at the link function scale) for the simulated data from both
266 distributions. We simulated a gradient of overdispersed data by adding noise to the linear
267 predictor with values from a Gaussian distribution with a mean of zero and ten standard
268 deviation values varying from 0 to 1. We evaluated test performance by comparing type I
269 error, power, and dispersion statistics across all parameter combinations in the simulated
270 datasets.

271 All models were fitted using the functions *glm* from the stats package or *glmer* from
272 the lme4 package (Bates et al., 2015) in R (v4.4; R Core Team, 2024). All dispersion tests
273 were performed with the DHARMA package (Hartig, 2026). For the simulation-based
274 residual variance test and the nonparametric Pearson residuals test, we set the number of
275 simulations to 250 (the default in DHARMA). All data, simulations and analysis codes are

276 available at the Zenodo repository (Leite, 2026). The supplementary material provides a
277 script file with instructions and examples for applying dispersion tests using the DHARMa
278 package.

279 *Theoretical expectations from simulations*

280 The classical (1) parametric Pearson residuals test assumes that the sample size (n -
281 asymptotic) and the expected values are sufficiently large (phi-asymptotic) (Venables &
282 Ripley, 2002). This implies that when the expected counts (or intercept) and/or the number
283 of observations are small, Pearson residuals may not provide reliable information about
284 model fit (see S2). Some corrections for Pearson residuals in small samples have been
285 suggested (e.g., Cordeiro, 2004; Cordeiro & Simas, 2009), but they are not currently
286 implemented in the most common R packages. Therefore, we expect the parametric
287 Pearson residuals test to perform well for GLMs, except in very small sample sizes and
288 expected counts (hereafter “small-data” situations).

289 It is unclear whether the parametric Pearson residuals test can be extended to
290 GLMMs or other hierarchical models, where counting residual degrees of freedom (df) is
291 not straightforward (Bolker et al., 2009; Luke, 2017). In mixed-effects models, the df
292 associated with a random effect are data-specific (adaptive shrinkage) and expected to lie
293 between one and the number of grouping levels (Baayen et al., 2008; Bolker et al., 2009;
294 Luke, 2017). Approaches exist to approximate df for random effects in LMMs (e.g.,
295 Schaalje et al., 2002), but their generalization to GLMMs remains an active area of
296 research. Current R packages that implement the parametric Pearson residuals test
297 approximate the df using the so-called naïve df (e.g., $n = 1$ per random effect) for testing

298 LMMs/GLMMs (Table 1). We expect the error introduced by this approximation to
299 increase with the number of random-effect groups. To test this, we varied the number of
300 groups in the random intercept (10, 50, and 100) in our simulated data.

301 In contrast to the parametric Pearson residuals test, we expect the (2) nonparametric
302 Pearson residuals test to be robust to small-data problems and to the presence of random
303 effects, as it doesn't rely on a specific parametric distribution. However, because the test
304 uses parametric bootstrapping, we expected it to run much more slowly than the other tests,
305 especially for more complex GLMMs. For this purpose, we compared the runtime of the
306 tests using a small set of simulated data (see S7).

307 For GLMMs tested with the (3) simulation-based residual variance test, we
308 compared the test's performance under the two simulation approaches, conditional and
309 unconditional on random effects. We expect lower power in the unconditional simulation
310 results, as overdispersion is a phenomenon at the model distribution level (i.e., at a higher
311 level). We evaluated the circumstances under which this test is reliable as a fast alternative
312 to both dispersion tests based on Pearson residuals.

313 *Case study 1: Redstart breeding pairs*

314 The first case study aims to show that empirical data analysis can give wrong
315 inferential results when dispersion is not taken into account. For that, we used the counts of
316 breeding pairs of the Common Redstart bird (*Phoenicurus phoenicurus*) in Switzerland
317 (Schmid et al. 1998), from the Swiss Breeding Bird Atlas 1993-1996, Swiss Ornithological
318 Institute. The data are openly available in the R package blmeco (Korner-Nievergelt et al.,
319 2019) as an accompanying dataset in Korner-Nievergelt et al. (2015). The data consist of

320 342 observations, and the counts are bird pairs within 1 km² plots. For our purposes, we
321 analyzed the effect of forest cover (%) on the abundance of breeding pairs, adding elevation
322 (in meters, with a quadratic term) as a covariate. We first applied the three dispersion tests
323 mentioned above to a Poisson GLM to check for dispersion problems and then fitted other
324 models accordingly.

325 *Case study 2: Wild and zoo-housed orangutan behavior*

326 The second case study analyzes more complex data and model situations in which
327 the parametric Pearson residuals test may fail and harm inference. The dataset and model
328 structures come from a study of behavioural differences in object exploration between wild
329 and zoo-housed Sumatran orangutans (Laumer et al., 2025). The original study evaluates
330 many aspects of object manipulation in wild and zoo-housed animals using complex
331 GLMMs (highly hierarchical and some with zero-truncated distributions). The authors
332 assessed only overdispersion, primarily using the parametric Pearson residuals test, and
333 found that for almost all models the dispersion parameters were below 1 (ranging from 0.33
334 to 0.78; Table S4b in the original study), yet they reported no dispersion issues. To compare
335 dispersion tests, we fitted and re-evaluated the model with the lowest dispersion parameter.
336 This model estimates the number of body parts involved in object manipulation using a
337 zero-truncated Poisson distribution, based on a dataset of 11,934 behavioral observations
338 from 445 day-animal observations across 51 animals. Predictors in these models included
339 the animal's age (with a quadratic effect) interacting with its origin (Zoo or Wild). Random
340 effects included random intercepts and random age slopes for individuals, as well as a day-
341 animal random intercept. We fitted the model using the *glmmTMB* R package (M. Brooks

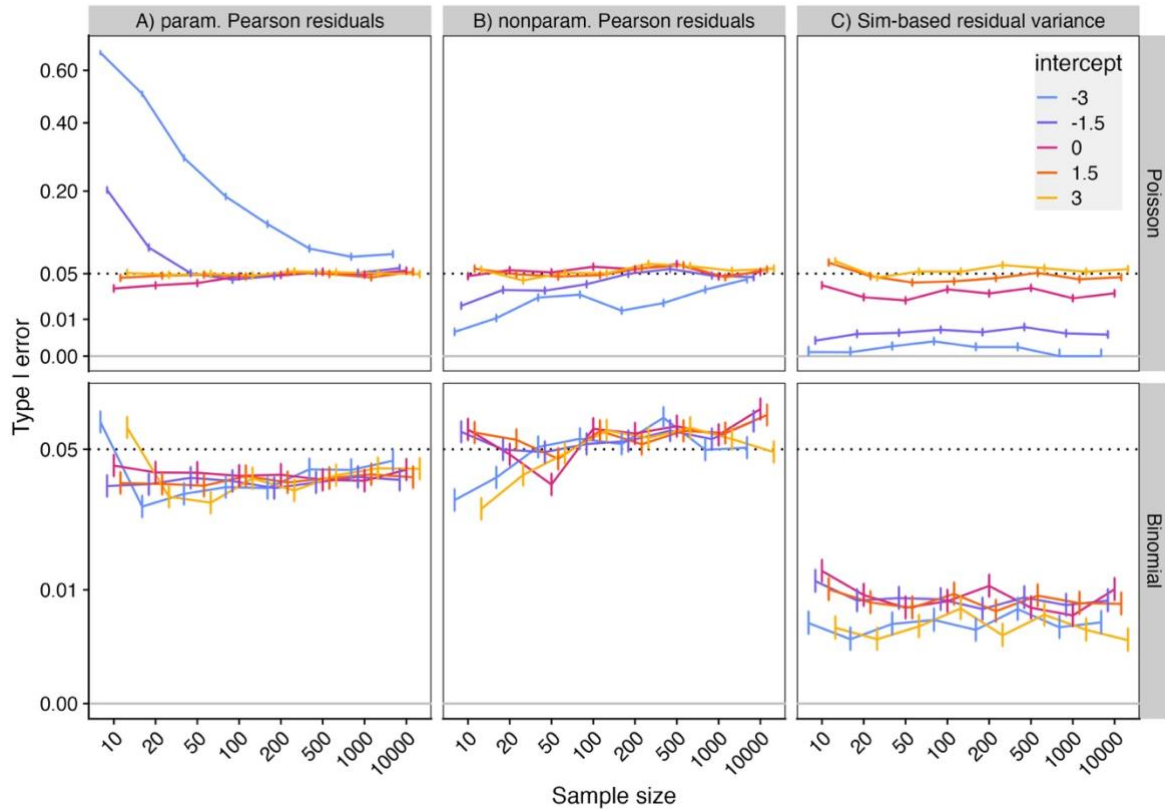
342 E. et al., 2017) and followed the model structure exactly as in the original paper (Laumer et
343 al., 2025). For the parametric Pearson residuals test, we fitted only for overdispersion
344 (alternative = “greater”), since the simulation results showed that this is the only possible
345 test for GLMMs (see results).

346 **Results**

347 *Performance on Poisson and binomial GLMs*

348 For Poisson GLMs, we found the expected distribution problems (Fig. 2): type I
349 error rates for the parametric Pearson residuals test were substantially high for the smallest
350 intercepts (-3), and they did not reach the nominal value of 0.05 even for very large sample
351 sizes ($n = 10,000$). The type I error rates for the nonparametric Pearson residuals test were
352 well calibrated, except for the smallest intercept (-3), with slightly conservative type I error
353 rates (< 0.05). For the simulation-based residual variance test, type I errors were
354 independent of sample size, but exhibited an intercept-dependent bias, ranging from almost
355 0 for the smallest intercept to 0.06 for the largest intercept.

356 For binomial GLMs, the type I error rates for the parametric Pearson residuals test
357 were generally conservatively calibrated around 0.04 (Fig. 2). Type I error rates for the
358 nonparametric Pearson residuals test averaged around 0.05 and 0.06, except for the very
359 low and very high intercepts (-3 and 3). For the simulation-based residual variance test,
360 type I error rates were conservatively very low for all simulated parameters, bouncing
361 below 0.01.



362

363 **Figure 2.** The simulation-based residual variance test has a more conservative type I error
 364 rate than both Pearson residual tests. The three dispersion tests were applied to Poisson
 365 (upper panels) and binomial N/K discrete proportion (lower panels, $K=10$ trials) GLMs:
 366 1A) parametric Pearson residuals test, B) nonparametric Pearson residuals test, and C)
 367 simulation-based residual variance test (see Table 1 for explanations). Simulations were run
 368 across different sample sizes (x-axis) and intercepts (colours, values on the link function
 369 scale). All points include a 95% confidence interval calculated from exact binomial tests
 370 across the 10,000 simulations. Note the square-root scale of the y-axis in the upper panel
 371 row. The dotted horizontal black line shows the 0.05 nominal type I error value.

372

The statistical power of the simulation-based residual variance test was lower than

373

the parametric and nonparametric Pearson residuals tests for both binomial and Poisson

374

GLMs, but tended to be similar with larger sample sizes (Fig. 3). We found that the reason

375

for this is the very conservative type I error rates (Fig. 2). When power is calibrated by

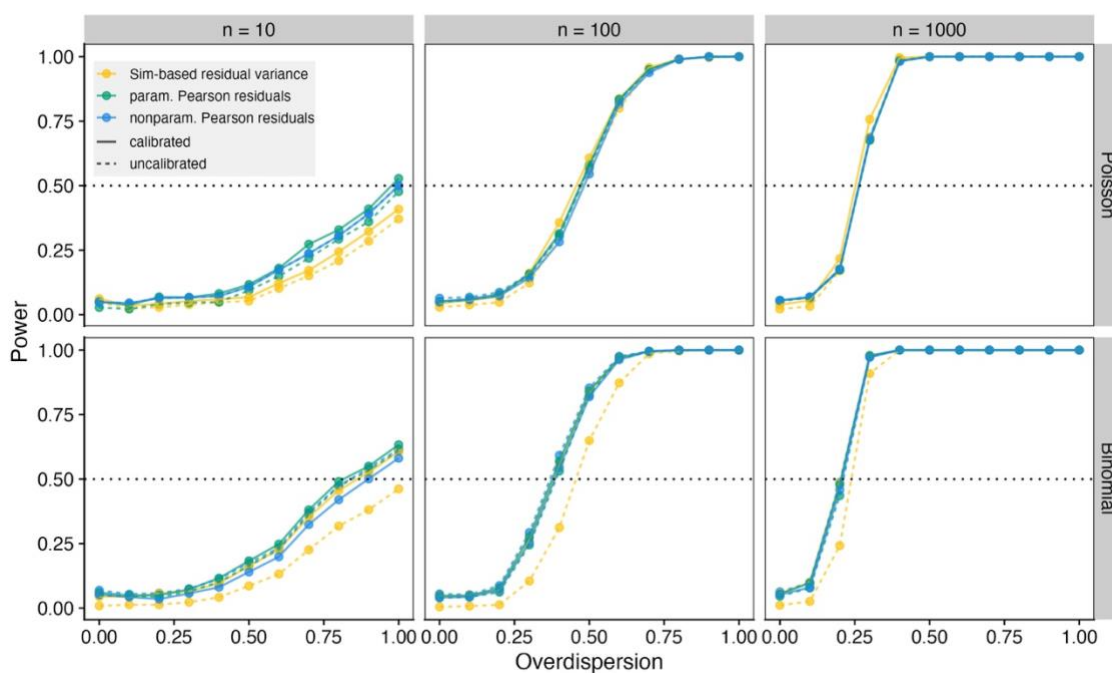
376

using the p-value at the 5% quantile of its empirical distribution for each simulation (details

377

in S6), the differences disappear (Fig. 3).

378 The dispersion statistics of the simulation-based residual variance test were highly
 379 dependent on the intercept, slope, and number of trials in the binomial model (see S5, Fig.
 380 S5.2), and tended to be smaller than those based on Pearson residuals. In contrast, for
 381 Poisson models, the values tended to be larger than those of Pearson statistics (Fig. S4.5).
 382 This may also explain the lower uncorrected power for the simulation-based residual
 383 variance test, especially for binomial models.

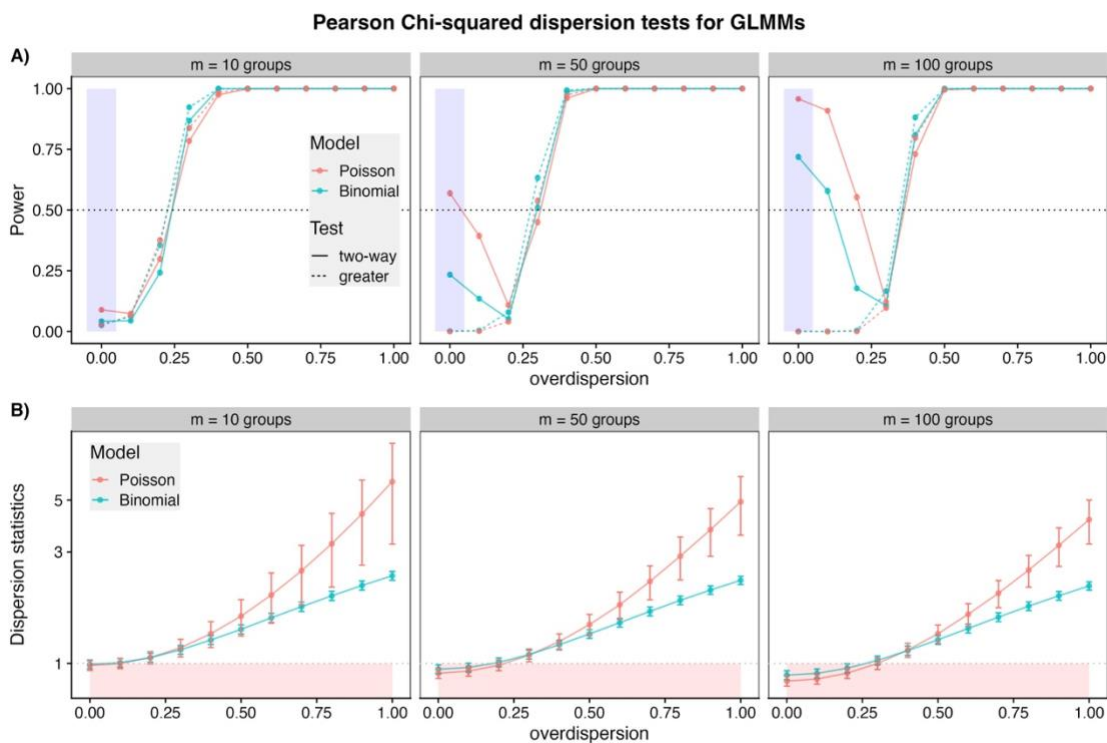


384
 385 **Figure 3.** The simulation-based residual variance test (in yellow) has lower power than
 386 both Pearson residuals tests (green and blue) for GLMs unless power is calibrated by type I
 387 error rates (dashed lines). Lower power is more evident for binomial models (upper panel)
 388 and smaller sample sizes (first two columns). Results based on 10,000 simulations per
 389 combination of parameters for an intercept = 0 and slope = 1. For all simulation results, see
 390 Fig. S6.1 and S6.2.

391 *GLMM performance*

392 For the GLMMs, we first compared the performance of the parametric Pearson
 393 residuals test (two-sided) for an increasing number of groups (m) in the random intercepts.
 394 As expected, the performance of the test failed for a large number of groups in the random

395 effects (Fig. 4A). The dispersion statistic was underestimated, and the type I error rates
 396 were too high because the test detected significant underdispersion. Testing only for
 397 overdispersion (“greater” test) when using the parametric Pearson residuals test appears to
 398 be the only reasonable approach for GLMMs (Fig. 4A). Still, it doesn’t prevent the
 399 dispersion statistics from being biased to lower values.

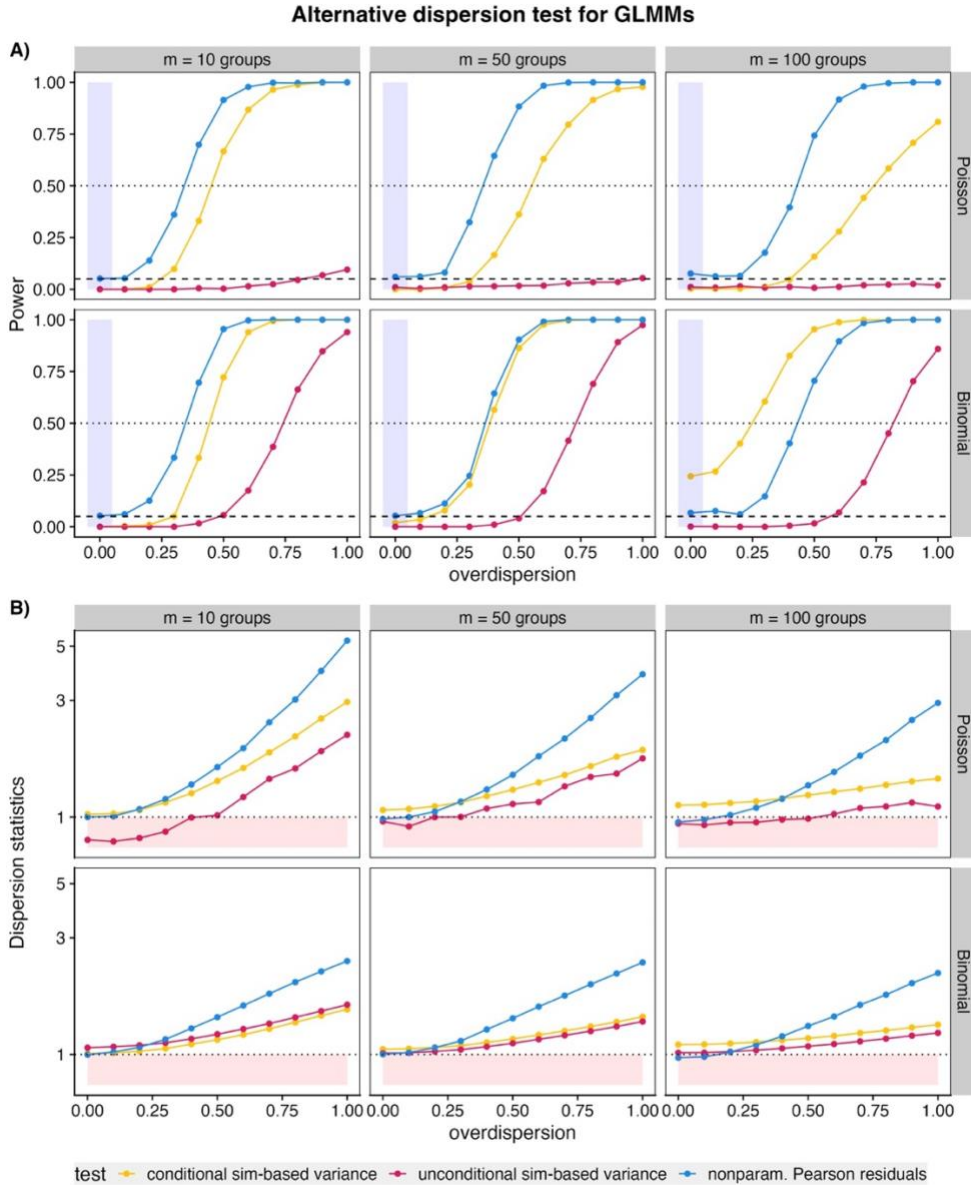


400

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

408 When comparing the alternative dispersion tests for GLMMs, the nonparametric
 409 Pearson residuals test showed very good results, with a type I error rate around 0.05 (Fig.
 410 S7.1 and S7.2) and higher power than the simulation-based residual variance tests (Fig. 5).

411 As expected, the unconditional simulation-based residual variance test had the worst
412 performance: very low type I errors (Fig. S7.1 and S7.2), very low power, and dispersion
413 statistics below 1 (Fig. 5B), especially for Poisson models. The conditional simulation-
414 based residual variance test also had very small type I errors (Fig. S7.1 and S7.2), but
415 power increased with the simulated overdispersion. The performance of both simulation-
416 based residual variance tests (unconditional and conditional) didn't change much with the
417 number of groups for the Poisson GLMMs, but it improved for the binomial GLMMs with
418 the increasing number of groups in the random intercept.



419

420 **Figure 5.** The nonparametric Pearson residuals test showed correct Type 1 error, higher
 421 power, and larger dispersion statistics than the simulated-based residual variance tests
 422 (conditional and unconditional to all random effects) for Poisson and binomial GLMMs.
 423 Power (A), type I error (shaded blue area in A), and dispersion statistics (B) for the
 424 alternative dispersion tests for Poisson and binomial GLMMs with different numbers of
 425 groups in random intercepts. The dashed horizontal line in (A) indicates the nominal value
 426 of 0.05 for type I error. The dotted horizontal line in (A) indicates the 50% power, and the
 427 dotted horizontal line in (B) indicates the dispersion statistics of 1. The results are based on
 428 1,000 simulations per parameter combination, with an intercept of 0 and a sample size (n)
 429 of 1,000.

430 *Case study 1: Redstart breeding pairs*

431 For case study 1, we initially found significant negative effects of elevation
432 (quadratic effect, $p=0.002$) and forest cover ($p=0.02$) on the number of breeding birds,
433 estimated using a Poisson GLM (Table S10.1). However, all three dispersion tests showed
434 similar dispersion parameters and p-values: significant overdispersion for the Poisson
435 GLM, with dispersion parameters ranging from 2.38 for both Pearson tests to 2.45 for the
436 simulation-based variance test (Table S10.2). In this case, we anticipate that many analysts
437 would switch to a negative binomial GLM, so we fitted such a model, and the
438 overdispersion was absent (nonsignificant dispersion parameters ranging from 1.02 to
439 1.03).

440 However, as explained in the introduction, overdispersion often arises through
441 misfit. In this case, a spatial autocorrelation check of the Poisson GLM revealed that the
442 true reason for the overdispersion was a spatial pattern in the residuals, which was also
443 present in the negative binomial GLM residuals (Table S10.3). By fitting a spatial
444 autocorrelated Poisson GLM, dispersion tests didn't detect any overdispersion as in the
445 negative binomial GLM. For both models, the slope for forest cover didn't change but
446 became nonsignificant (Table S10.1). In this scenario, we would prefer to interpret the
447 spatial Poisson GLM rather than changing the distribution.

448 *Case study 2: Wild and zoo-housed orangutan behavior*

449 In case study 2, we demonstrate that the parametric Pearson residuals tests were
450 inadequate to assess dispersion problems in the evaluated complex GLMM model. The
451 zero-truncated Poisson model showed residual underdispersion under both the simulation-
452 based and nonparametric Pearson tests. A reassessment of this model with a zero-truncated

453 Conway-Maxwell-Poisson distribution showed no dispersion problems according to the
 454 simulation-based residual variance test. For this distribution in *glmmTMB*, since the
 455 Pearson residuals are not computed, Pearson residual tests are not applicable. Although the
 456 zero-truncated Conway-Maxwell-Poisson GLMM fits much better to the data (e.g., AIC
 457 difference of 844), the ecological inference (significance of animals' origin and age) and
 458 qualitative results remained similar to the original zero-truncated Poisson GLMM.

459 **Table 2.** Dispersion tests applied to the orangutan data analysis from Laumer et al. (2025).
 460 The original model estimated the number of body parts involved in object manipulation
 461 using a zero-truncated Poisson GLMM. We compared the dispersion tests applied to this
 462 model, since the original parametric Pearson residuals test cannot assess underdispersion.
 463 After detecting underdispersion in the previous model, we fitted a zero-truncated Conway-
 464 Maxwell-Poisson (CMP) GLMM and retested for dispersion. However, for the zero-
 465 truncated CMP, the only option is the simulation-based residual variance test, as the
 466 Pearson residuals are not available.

Model	Dispersion test	Parameter	P.value
Zero-truncated Poisson	Sim-based res. variance (two-sided)	0.708	0.00
	Parametric Pearson res. (greater)	0.330	1.00
	Nonparametric Pearson res. (two-sided)	0.680	0.00
Zero-truncated CMP	Sim-based res. variance (two-sided)	1.021	0.15

467 Discussion

468 Through extensive simulation analysis, we identified the strengths and limitations of
 469 three dispersion tests that, in principle, could be widely applicable across GLM and GLMM
 470 distributions commonly used in ecological data analysis. We conclude that the
 471 nonparametric Pearson residuals test is the most reliable general test currently available.
 472 For GLMs, this test exhibited similar power to the parametric Pearson residuals test but
 473 with more reliable type I error rates in small-sample situations. The downside of this test is
 474 that it can be computationally expensive, with runtimes in the order of minutes for larger
 475 GLMMs, and, as our case study 2 showed, Pearson residuals may not be available for non-

476 standard distributions even in commonly used R packages. The simulation-based residual
477 variance test for GLMs is more generally applicable and fast to compute, but its dispersion
478 statistic is more difficult to interpret and often leads to overly conservative type I errors.
479 This results in low power unless it is additionally calibrated using a simulated p-value
480 distribution. The parametric Pearson residuals test is computationally efficient, but it is
481 unreliable in small-data situations and in the presence of random effects. Below, we discuss
482 these points in more detail and, based on our case studies, provide recommendations for
483 general users who rely on already-implemented R packages for model fitting and
484 diagnostics.

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

486 We showed that the parametric Pearson residuals test, although popular, quick, and
487 relatively easy to compute, has two main disadvantages: it performs poorly in (1) small-
488 data situations (Fig. 2) and (2) in the presence of random effects (Fig. 4). The first problem
489 arises from a mismatch between the distribution of the Pearson statistic and the chi-squared
490 distribution under small-data conditions (Fig. S2.1 and S2.2). This phenomenon has already
491 been studied (e.g., Fletcher, 2012; Kuss, 2002), with suggested corrections (Farrington,
492 1996; McCullagh, 1985). However, none of these corrections are implemented in the
493 current R packages (Table 1), and we believe it will be difficult to devise corrections that
494 work across a wide range of distributions.

495 The second problem arises because counting 1 degree of freedom (df) for a random
496 effect, as done in most implementations of this test, typically underestimates the true model
497 df , and this underestimation increases with the number of levels of the random effect. The

498 result is a bias in the dispersion statistic towards underdispersion that increases with the
499 number of random-effect levels (Fig. 4). Two-sided tests would therefore often wrongly
500 detect significant underdispersion in perfectly valid GLMMs, which is likely why most R
501 implementations of this test only test for overdispersion. When applying this test to
502 GLMMs, we recommend following the same approach and ignoring dispersion statistics
503 smaller than 1. Nevertheless, this is an unsatisfactory solution, as the biased dispersion
504 statistic also causes a loss of power.

505 A possible solution for GLMMs could be to use a better approximation of the
506 residual degrees of freedom (df). For LMMs, approximations for denominator df have been
507 successfully used for hypothesis testing (Luke, 2017), for example, the Satterthwaite (1946)
508 and the Kenward-Roger (2009). Although there is some evidence that these approximations
509 are also accurate for GLMMs (Stroup, 2015), the main R packages implementing these
510 methods are currently limited to LMMs (e.g., *pbkrtest* Halekoh & Højsgaard, 2014;
511 *lmerTest* Kuznetsova et al., 2017). However, the recently released package *glmmrBase*
512 (Watson, 2024) allows these methods to be applied to GLMMs. We performed some
513 parametric Pearson residuals tests for Poisson GLMMs using a modified residual df
514 approximation (see S9). Although the parametric Pearson residuals tests with the
515 approximated residual df performed much better than those with the naïve residual df , they
516 still underperformed compared to the nonparametric Pearson residuals test when there were
517 a large number of groups in the random effects (Fig. S9.4), especially in small-data
518 situations.

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

520 The simulation-based residual variance test developed in the R package DHARMA
521 (Hartig, 2026) is the main alternative to the family of Pearson residuals tests. Its principle is
522 simple: when the model is correctly specified, the variance of the observed data should
523 match that of data simulated from the model. The main advantage of this approach is that it
524 is a nonparametric test applicable to any model structure and does not require refitting the
525 model, making it considerably faster and easier to implement. We also note that for
526 GLMMs, simulations should be performed conditionally to avoid a loss of power,
527 presumably due to the increased variability created by re-simulating the random effects
528 (unconditional simulations).

529 The disadvantages of this approach are that it is often overly conservative, resulting
530 in lower power than the Pearson residuals tests. Additionally, the calculated dispersion
531 statistic differs from the Pearson dispersion statistic, making it difficult to compare the two
532 approaches. We conjectured that both problems could be related to the test statistic being
533 based on the raw variance (rather than a scaled variance, as with the Pearson statistics),
534 which may overrepresent observations with large values. We considered scaling each
535 observation by the expected variance, but this is not readily available for a wide class of
536 models, and using simulations to approximate it failed for discrete-valued distributions (see
537 S8).

538 **Conclusions and recommendations**

539 Although neither of the options considered for testing dispersion excelled in all
540 dimensions (Fig. 6), our primary recommendation is that, for standard GLMs with
541 sufficient data, the parametric Pearson chi-squared test, available in many packages (Table

542 1), can be safely used, as shown in case study 1. In complex situations, particularly for
 543 GLMMs, we recommend the nonparametric Pearson residuals test. It has very few
 544 weaknesses, other than being computationally costly. If the nonparametric Pearson
 545 residuals test cannot be calculated due to speed or convergence problems with refitting
 546 complex models, we recommend using the simulation-based residual variance test with
 547 simulations performed conditionally on the fitted random effects. All three approaches are
 548 available via the *testDispersion* function in the DHARMA R package (Hartig, 2026). We
 549 provide a supplementary file with instructions and an example for applying dispersion tests
 550 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	+++	-	+	-	+++

551
 552 **Figure 6.** Performance comparisons of the dispersion tests evaluated for “dimension” of
 553 Poisson and binomial models: GLMs in general, GLMs with small sample size or intercept
 554 (“small data”), GLMMs with one random effect with few groups/levels, GLMMs with
 555 many groups/levels in a random effect, and computational time for calculating the test
 556 (speed). The symbols mean: “-” bad performance, “+” good performance, “+++” very good
 557 performance.

558 Although our simulation examples focused on overdispersion, the tests considered
 559 in our study can also be used to detect underdispersion by testing the dispersion “two-
 560 sided” or “less than” against null statistics. We show this usage in case study 2, where tests

561 detected underdispersion. The clear exception is testing for underdispersion using the
562 parametric Pearson residuals test for GLMMs, which would be anti-conservative due to the
563 discussed bias towards underdispersion in the presence of random effects.

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

565 For interpretation and applied ecological data analysis, we stress that a significant
566 over- or underdispersion result does not necessarily indicate that the distribution must be
567 changed. First, hypothesis tests evaluate statistical rather than ecological significance. In
568 other words, a significant test for overdispersion indicates that the overdispersion signal
569 deviates from a null expectation, but the p-value does not measure the strength of the
570 deviation. The first step in a dispersion test should thus be to examine how much the
571 dispersion statistic deviates from the expected value of 1. For very large sample sizes, small
572 departures from 1 may be statistically significant, but they may not necessarily warrant a
573 change to the model. Second, after finding that a dispersion problem is both significant and
574 meaningful, we suggest checking for problems beyond the distribution, such as
575 spatial/temporal/phylogenetic autocorrelation, heteroscedasticity, missing predictors, an
576 incorrect link function, excess zeros, or overfitting. As we found in case study 1 on spatial
577 autocorrelation, these types of model misspecifications often cause over-/underdispersion,
578 but can be distinguished from a “real” distributional problem through careful residual
579 checks. Blindly changing the distribution only masks the problem, without offering a real
580 solution to the underlying problems.

581 Finally, after ruling out potential model misspecifications leading to under-
582 /overdispersion, we should consider changing the model’s distribution, as we may be facing

583 an ‘intrinsic’ under-/overdispersion problem, likely due to the nature of ecological data
584 (Box 1). A traditional and flexible solution is to use the ‘quasi’ distributions (Wedderburn,
585 1974), which essentially correct p-values but do not represent an explicit data-generating
586 process with an associated likelihood, precluding, for example, simulation from the fitted
587 model. A second alternative for adding dispersion is to use observation-level random
588 effects (Bolker et al., 2009; Elston et al., 2001; Harrison, 2014; Ozgul et al., 2009). While
589 often a reasonable solution, excessive use of random effects can create problems in
590 calculating other statistical indicators (such as p-values) that we would rather avoid. For
591 that reason, we consider the best solution to address ‘intrinsic’ under-/overdispersion is to
592 switch to the corresponding variable-dispersion distributions. For overdispersed count data,
593 the most used is the negative binomial (see S1). However, other distributions have been
594 used in ecology to handle both over- and underdispersion, such as the generalized Poisson,
595 the Conway-Maxwell-Poisson, the Double Poisson, and the Good distributions (Agis et al.,
596 2024; M. E. Brooks et al., 2019; Lynch et al., 2014). For discrete proportions data, the beta-
597 binomial distribution (Harrison, 2015) is considered the most appropriate for overdispersed
598 binomial models (Harrison, 2015). Regardless of the approach, an “over-/underdispersion-
599 free” GLM/GLMM is essential for better interpretation of ecological models and for
600 facilitating sound scientific discoveries.

601 **Acknowledgements**

602 This study was funded by the Deutsche Forschungsgemeinschaft (DFG), project number
603 528747641.

604 **Author contributions**

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

609 **Conflict of interest**

610 The authors, MSL and FH, are developers of the R package DHARMa, which implements
611 the dispersion tests used in this study.

612 **References**

- 613 Agis, D., Tur, J., Moriña, D., Puig, P., & Fernández-Fontelo, A. (2024). good: An R
614 package for modelling count data. *Methods in Ecology and Evolution*, *15*(12),
615 2192–2197. <https://doi.org/10.1111/2041-210X.14387>
- 616 Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with
617 crossed random effects for subjects and items. *Journal of Memory and Language*,
618 *Special Issue: Emerging Data Analysis*, *59*(4), 390–412.
619 <https://doi.org/10.1016/j.jml.2007.12.005>
- 620 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for
621 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and*
622 *Language*, *68*(3), 255–278. <https://doi.org/10.1016/j.jml.2012.11.001>

- 623 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects
624 Models Using lme4. *Journal of Statistical Software*, 67(1), 1–48.
625 <https://doi.org/10.18637/jss.v067.i01>
- 626 Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H. H.,
627 & White, J.-S. S. (2009). Generalized linear mixed models: A practical guide for
628 ecology and evolution. *Trends in Ecology & Evolution*, 24(3), 127–135.
629 <https://doi.org/10.1016/j.tree.2008.10.008>
- 630 Brooks, M., E., Kristensen, K., Benthem, K., J., van, Magnusson, A., Berg, C., W., Nielsen,
631 A., Skaug, H., J., Mächler, M., & Bolker, B., M. (2017). glmmTMB Balances Speed
632 and Flexibility Among Packages for Zero-inflated Generalized Linear Mixed
633 Modeling. *The R Journal*, 9(2), 378. <https://doi.org/10.32614/RJ-2017-066>
- 634 Brooks, M. E., Kristensen, K., Darrigo, M. R., Rubim, P., Uriarte, M., Bruna, E., & Bolker,
635 B. M. (2019). Statistical modeling of patterns in annual reproductive rates. *Ecology*,
636 100(7), e02706. <https://doi.org/10.1002/ecy.2706>
- 637 Cameron, A. C., & Trivedi, P. (2023). *overdisp: Overdispersion in count data multiple*
638 *regression analysis* (Version 0.1.2) [Computer software].
639 <https://doi.org/10.32614/CRAN.package.overdisp>
- 640 Cameron, A. C., & Trivedi, P. K. (1990). Regression-based tests for overdispersion in the
641 Poisson model. *Journal of Econometrics*, 46(3), 347–364.
642 [https://doi.org/10.1016/0304-4076\(90\)90014-K](https://doi.org/10.1016/0304-4076(90)90014-K)

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

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

647 Collings, B. J., & Margolin, B. H. (1985). Testing goodness of fit for the Poisson
648 assumption when observations are not identically distributed. *Journal of the*
649 *American Statistical Association*, 80(390), 411–418.

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

652 Cordeiro, G. M., & Simas, A. B. (2009). The distribution of Pearson residuals in
653 generalized linear models. *Computational Statistics & Data Analysis*, 53(9), 3397–
654 3411. <https://doi.org/10.1016/j.csda.2009.02.025>

655 Dean, C. (1992). Testing for overdispersion in Poisson and binomial regression models.
656 *Journal of the American Statistical Association*, 87(418), 451–457.
657 <https://doi.org/10.2307/2290276>

658 Dean, C., & Lawless, J. F. (1989). Tests for detecting overdispersion in Poisson regression
659 models. *Journal of the American Statistical Association*, 84(406), 467–472.
660 <https://doi.org/10.1080/01621459.1989.10478792>

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

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

666 Farrington, C. P. (1996). On assessing goodness of fit of generalized linear models to sparse
667 data. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(2),
668 349–360.

669 Feng, C., Li, L., & Sadeghpour, A. (2020). A comparison of residual diagnosis tools for
670 diagnosing regression models for count data. *BMC Medical Research Methodology*,
671 20(1), 175. <https://doi.org/10.1186/s12874-020-01055-2>

672 Fisher, R. A. (1950). The significance of deviations from expectation in a Poisson series.
673 *Biometrics*, 6(1), 17–24. <https://doi.org/10.2307/3001420>

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

676 Gómez-Rubio, V., Ferrándiz-Ferragud, J., & López-Quílez, A. (2005). Detecting clusters of
677 disease with R. *Journal of Geographical Systems*, 7(2), 189–206.
678 <https://doi.org/10.1007/s10109-005-0156-5>

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

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

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

688 Hartig, F. (2026). *DHARMa: Residual Diagnostics for Hierarchical (Multi-Level / Mixed)*
689 *Regression Models* (Version 0.5.0) [Computer software]. [https://CRAN.R-](https://CRAN.R-project.org/package=DHARMa)
690 [project.org/package=DHARMa](https://CRAN.R-project.org/package=DHARMa)

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

694 Hilbe, J. M. (2011). *Negative Binomial Regression* (2nd ed.). Cambridge University Press.

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

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

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

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

705 Kleiber, C., & Zeileis, A. (2008). *Applied econometrics with R [R package AER version*
706 *1.2-14]*. Springer-Verlag. <https://doi.org/10.1007/978-0-387-77318-6>

- 707 Korner-Nievergelt, F., Roth, T., Felten, S. von, Guelat, J., Almasi, B., & Korner-Nievergelt,
708 P. (2019). *blmeco: Data Files and Functions Accompanying the Book “Bayesian*
709 *Data Analysis in Ecology using R, BUGS and Stan”* (Version 1.4) [Computer
710 software]. <https://cran.r-project.org/web/packages/blmeco/index.html>
- 711 Korner-Nievergelt, F., Von Felten, S., Roth, T., Almasi, B., Guélat, J., & Korner-Nievergelt,
712 P. (2015). *Bayesian data analysis in ecology using linear models with R, BUGS, and*
713 *Stan*. Elsevier/AP, Academic Press is an imprint of Elsevier.
- 714 Kuss, O. (2002). Global goodness-of-fit tests in logistic regression with sparse data.
715 *Statistics in Medicine, 21*(24), 3789–3801. <https://doi.org/10.1002/sim.1421>
- 716 Kuznetsova, A., Brockhoff, P., & Christensen, R. (2017). lmerTest package: Tests in linear
717 mixed effects models. *Journal of Statistical Software, Articles, 82*(13).
718 <https://doi.org/10.18637/JSS.V082.I13>
- 719 Lai, J., Lortie, C. J., Muenchen, R. A., Yang, J., & Ma, K. (2019). Evaluating the popularity
720 of R in ecology. *Ecosphere, 10*(1), e02567. <https://doi.org/10.1002/ecs2.2567>
- 721 Laumer, I. B., Kansal, S., Van Cauwenberghe, A., Rahmaeti, T., Setia, T. M., Mundry, R.,
722 Haun, D., & Schuppli, C. (2025). Wild and zoo-housed orangutans differ in how
723 they explore objects. *Scientific Reports, 15*(1), 14853.
724 <https://doi.org/10.1038/s41598-025-97926-z>
- 725 Lawless, J. F. (1987). Negative binomial and mixed Poisson regression. *Canadian Journal*
726 *of Statistics, 15*(3), 209–225. <https://doi.org/10.2307/3314912>

727 Leite, M. de S. (2026). *Data, simulations and code from: Dispersion tests in generalized*
728 *linear mixed-effects models - a methods comparison and practical guide for*
729 *ecologists* [Computer software]. Zenodo. <https://doi.org/10.5281/zenodo.17611060>

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

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

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

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

742 Lynch, H. J., Thorson, J. T., & Shelton, A. O. (2014). Dealing with under- and over-
743 dispersed count data in life history, spatial, and community ecology. *Ecology*,
744 *95*(11), 3173–3180. <https://doi.org/10.1890/13-1912.1>

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

- 748 McCullagh, P., & Nelder, J. (1989). Generalized linear models. *Journal of the Royal*
749 *Statistical Society*, 135(3), 370–384.
- 750 McMahon, S. M., & Diez, J. M. (2007). Scales of association: Hierarchical linear models
751 and the measurement of ecological systems. *Ecology Letters*, 10(6), 437–452.
752 <https://doi.org/10.1111/j.1461-0248.2007.01036.x>
- 753 Moral, R. A., Hinde, J., & Demétrio, C. G. B. (2017). Half-normal plots and overdispersed
754 models in R: The hnp package. *Journal of Statistical Software*, 81, 1–23.
755 <https://doi.org/10.18637/jss.v081.i10>
- 756 Nakagawa, S., Ortega, S., Gazzea, E., Lagisz, M., Lenz, A., Lundgren, E., & Mizuno, A.
757 (2026). Location–scale models in ecology and evolution: Heteroscedasticity in
758 continuous, count and proportion data. *Methods in Ecology and Evolution*, 17(2),
759 554–566. <https://doi.org/10.1111/2041-210x.70203>
- 760 Ohara Hines, R. J. (1997). A comparison of tests for overdispersion in generalized linear
761 models. *Journal of Statistical Computation and Simulation*, 58(4), 323–342.
762 <https://doi.org/10.1080/00949659708811838>
- 763 Ozgul, A., Oli, M. K., Bolker, B. M., & Perez-Heydrich, C. (2009). Upper respiratory tract
764 disease, force of infection, and effects on survival of gopher tortoises. *Ecological*
765 *Applications*, 19(3), 786–798.
- 766 Papadakis, M., Tsagris, M., Fafalios, S., Dimitriadis, M., & Lasithiotakis, M. (2025).
767 *Rfast2: A collection of efficient and extremely fast R functions II* (Version 0.1.5.4)
768 [Computer software]. <https://CRAN.R-project.org/package=Rfast2>

769 Puig, P., Valero, J., & Fernández-Fontelo, A. (2024). Some mechanisms leading to
770 underdispersion: Old and new proposals. *Scandinavian Journal of Statistics*, 51(1),
771 245–267. <https://doi.org/10.1111/sjos.12677>

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

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

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

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

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

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

791 Stroup, W. W. (2015). Rethinking the analysis of non-normal data in plant and soil science.
792 *Agronomy Journal*, 107(2), 811–827. <https://doi.org/10.2134/agronj2013.0342>

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

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

798 Watson, S. I. (2024). *Generalised linear mixed model specification, analysis, fitting, and*
799 *optimal design in R with the glmmr packages* (arXiv: 2303.12657). arXiv.
800 <https://doi.org/10.48550/arXiv.2303.12657>

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

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

806 Yang, Z., Hardin, J. W., Addy, C. L., & Vuong, Q. H. (2007). Testing approaches for
807 overdispersion in Poisson regression versus the generalized Poisson model.
808 *Biometrical Journal*, 49(4), 565–584. <https://doi.org/10.1002/bimj.200610340>

809 Zeileis, A., & Hothorn, T. (2002). Diagnostic checking in regression relationships. *R News*,
810 2(3), 7–10.

811

1 **Supporting information for:**

2 **Dispersion tests in generalized linear mixed-effects models: a**

3 **methods comparison and practical guide for ecologists**

4

5 **S1. Trend analysis and current ecological literature practices on**

6 **dispersal issues**

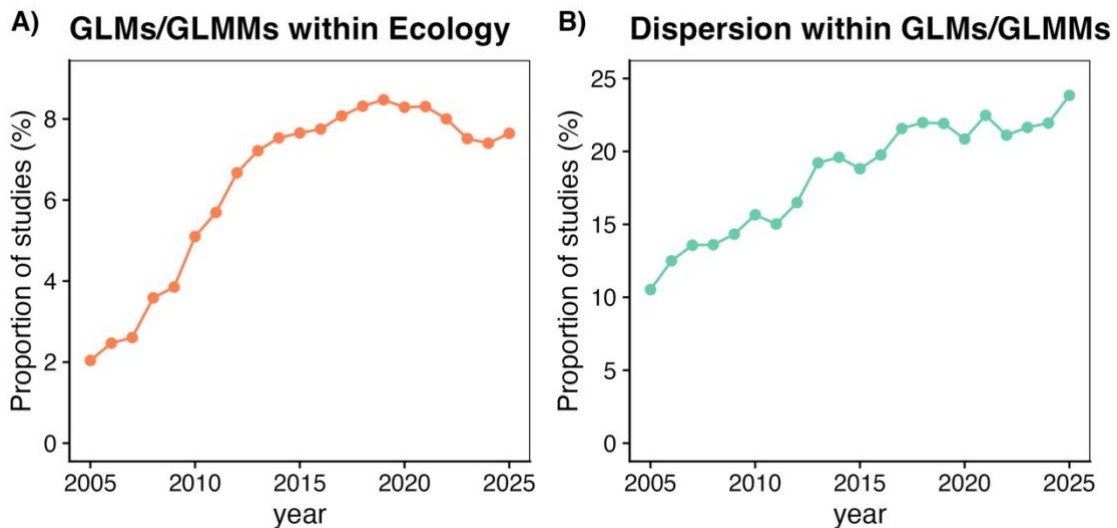
7 To understand the extent of ecological studies relying on GLMs/GLMMs for count and
8 discrete proportion data and those that address dispersion issues, we conducted a text analysis of
9 the ecological literature over the past 20 years. We used the R package ‘europepmc’ (v0.4.3,
10 Jahn, 2023) to search for articles in the PubMed and Medline NLM databases from 2005 to
11 2025. We used combinations of words (Table S1.1) to retrieved the annual records for: (1)
12 the percentage of ecological papers using GLMs/GLMMs for count and discrete
13 proportion data (Figure S1.1a), (2) the percentage of those papers that mention
14 dispersion terms in general (Figure S1.1b), (3) the percentage of ecological papers using
15 GLMs/GLMMs for count data mentioning dispersion terms (Figure BOX 1, main text),
16 and (4) the percentage of ecological papers using GLMs/GLMMs for discrete
17 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.



27

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

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

44 Of the 40 papers explicitly testing dispersion, 25 reported using the DHARMA R
 45 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 generalized 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.

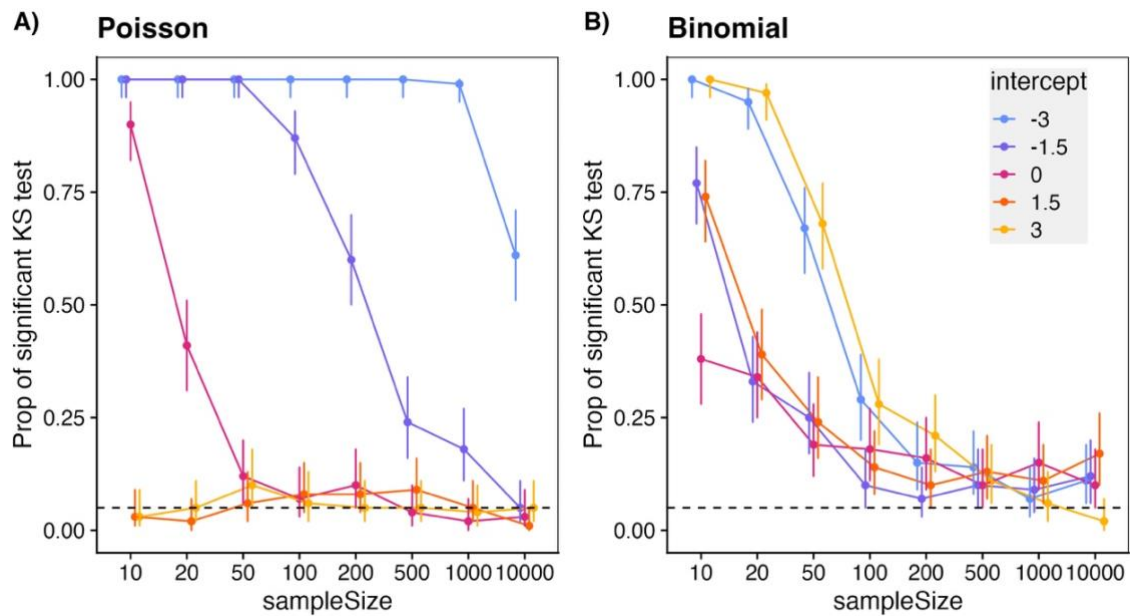
64 **S2. Pearson statistics and Chi-squared distribution**

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

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

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

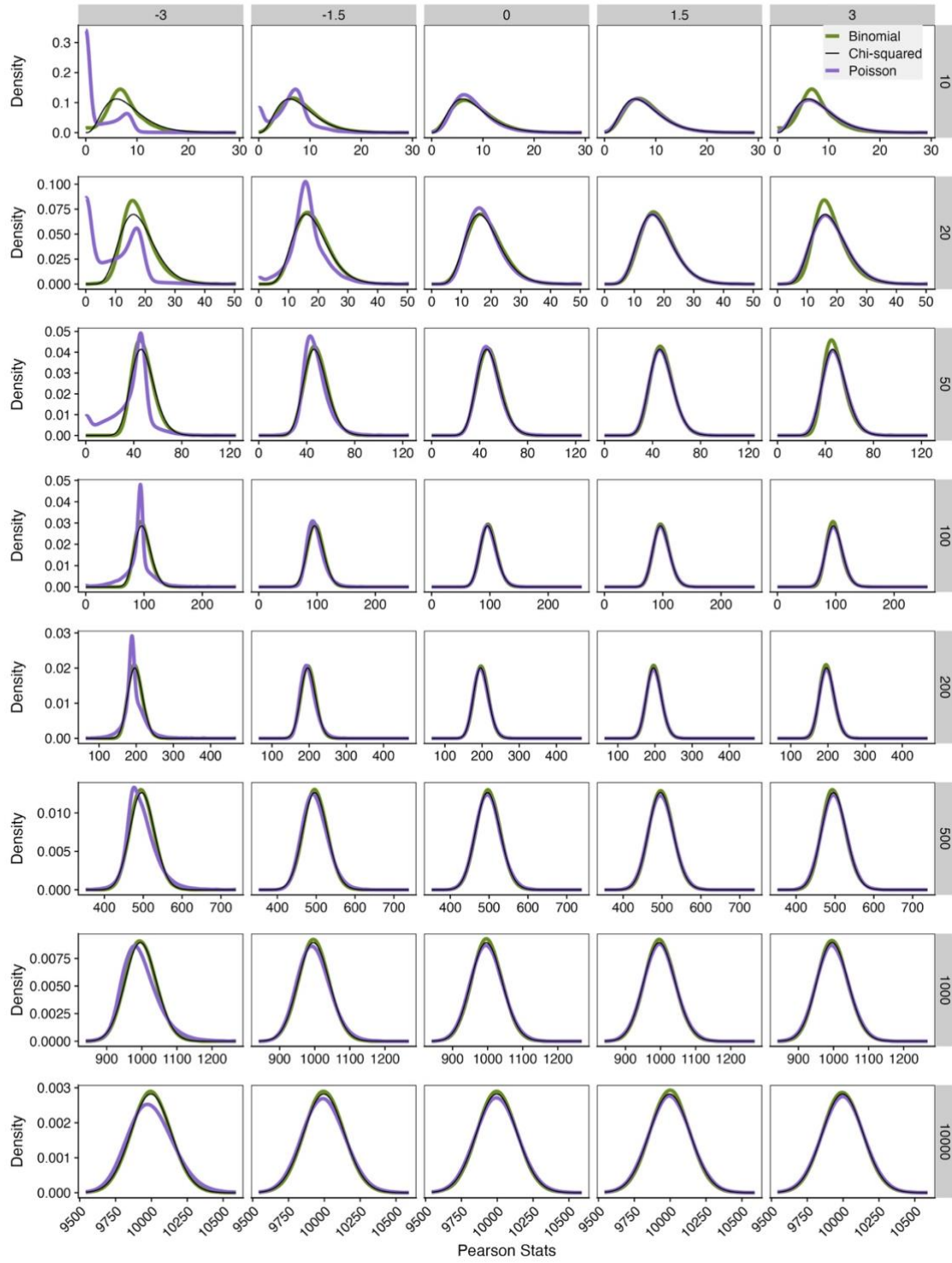
88 sample sizes (10, 20, 50) (Figure S2.1B). The proportion of significant KS tests
 89 decreased with sample size, but did not reach the nominal value of 0.05, even for very
 90 large sample sizes and intermediate intercept values (-1.5, 0, 1.5).



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



99

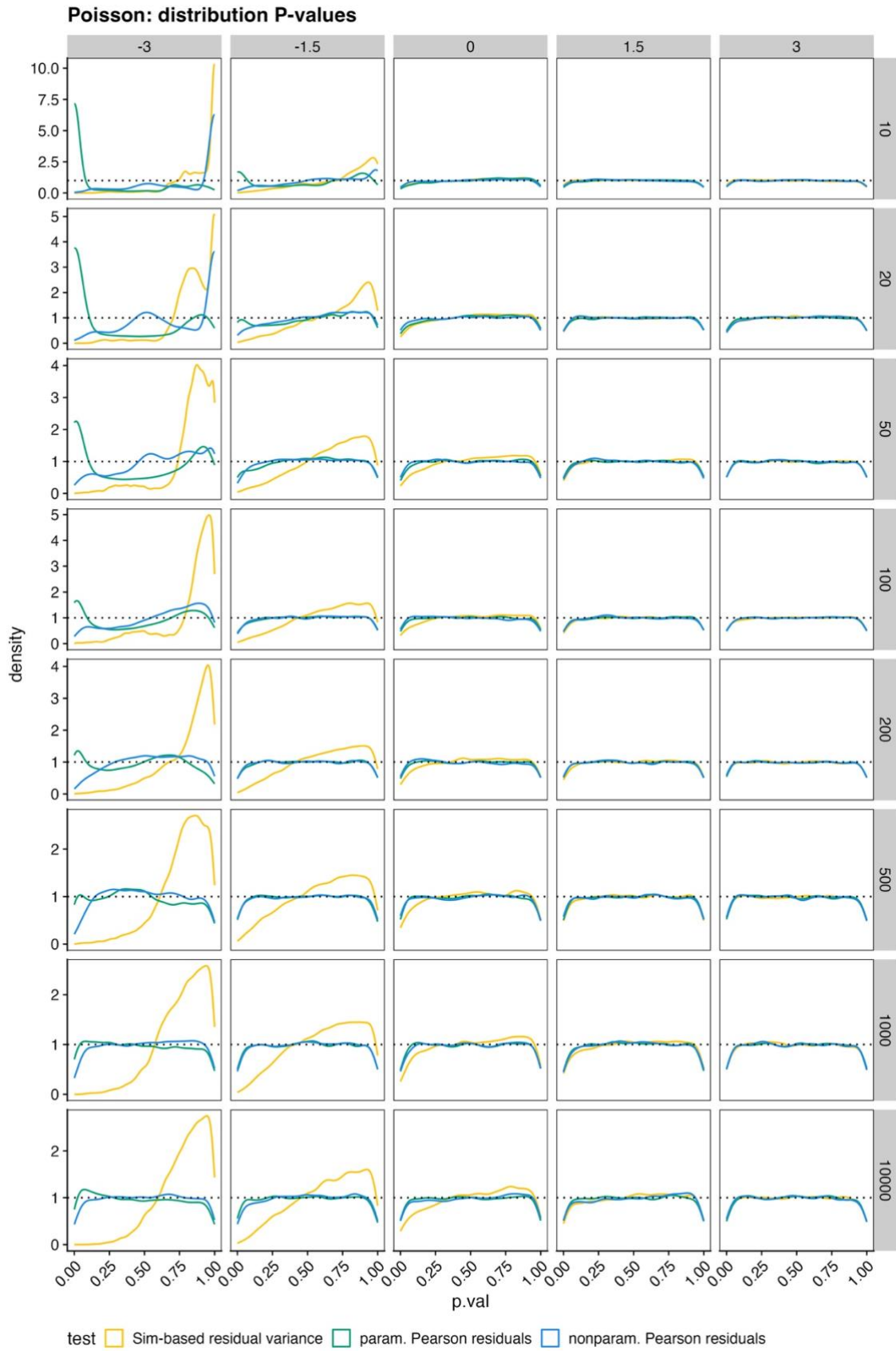
100 **Figure S2.2.** Mean Pearson statistics distribution (from 100 simulated curves) for the
 101 binomial (green) and Poisson (purple), and the Chi-square distribution in black.

102 **S3. Type I error rates for the GLMs**

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

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

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

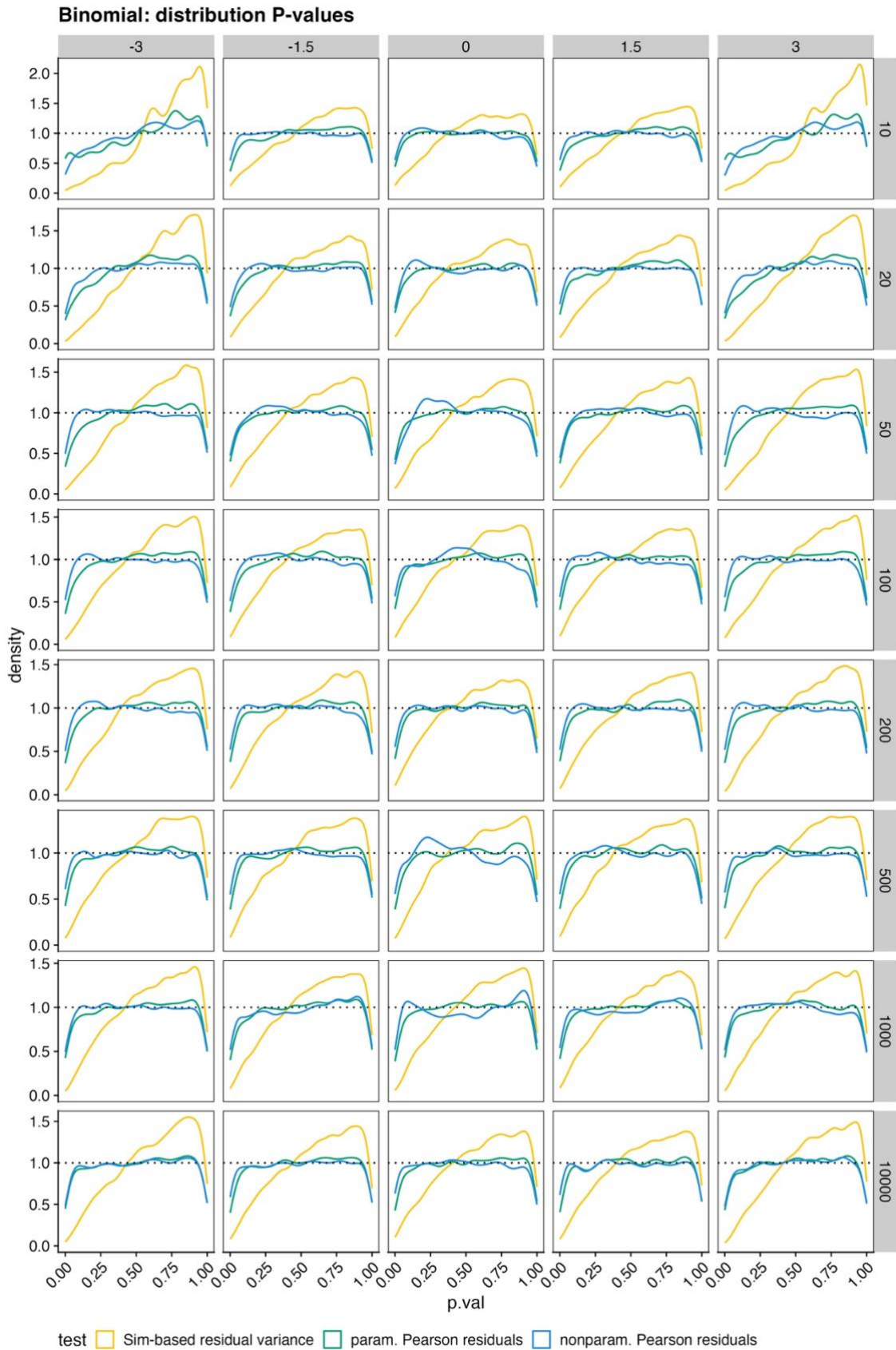


123

124

125

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).

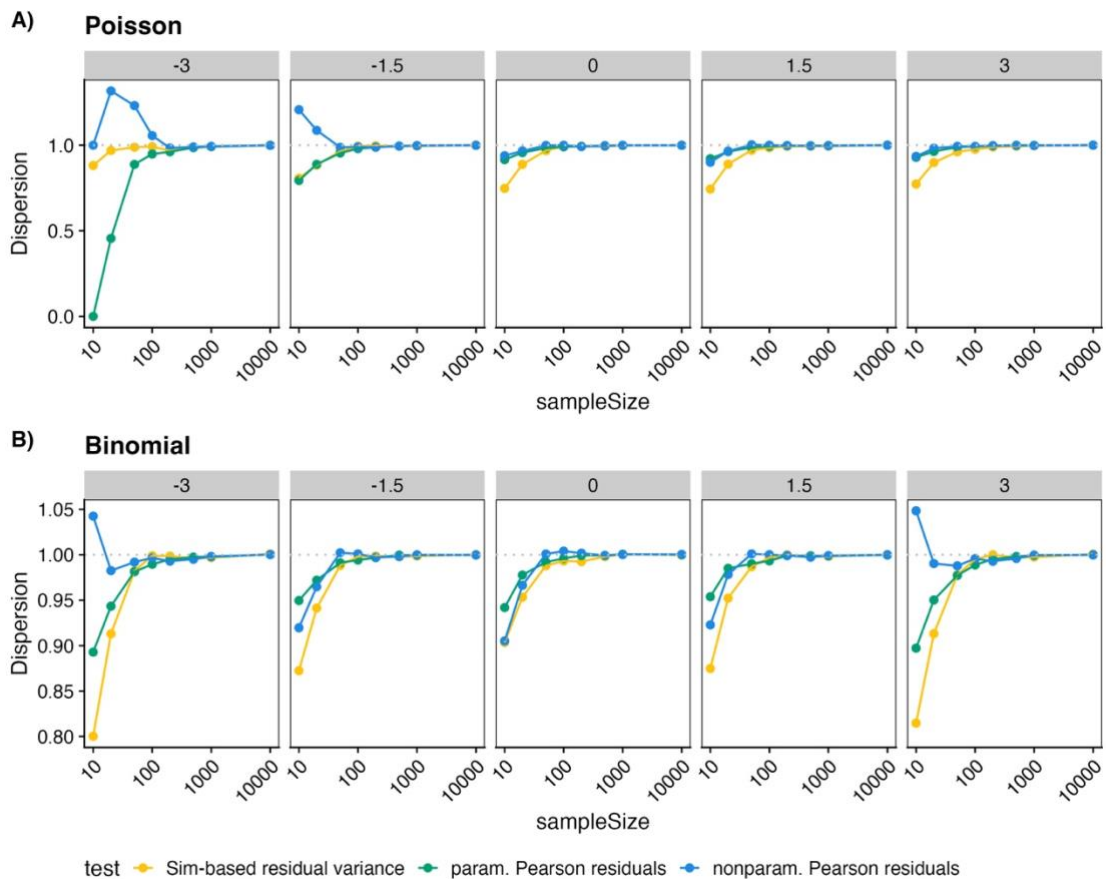


126

127 **Figure S3.2.** Distribution of p-values for the binomial GLMs for each dispersion test.
 128 10,000 simulations per simulation set (intercept x sample size).

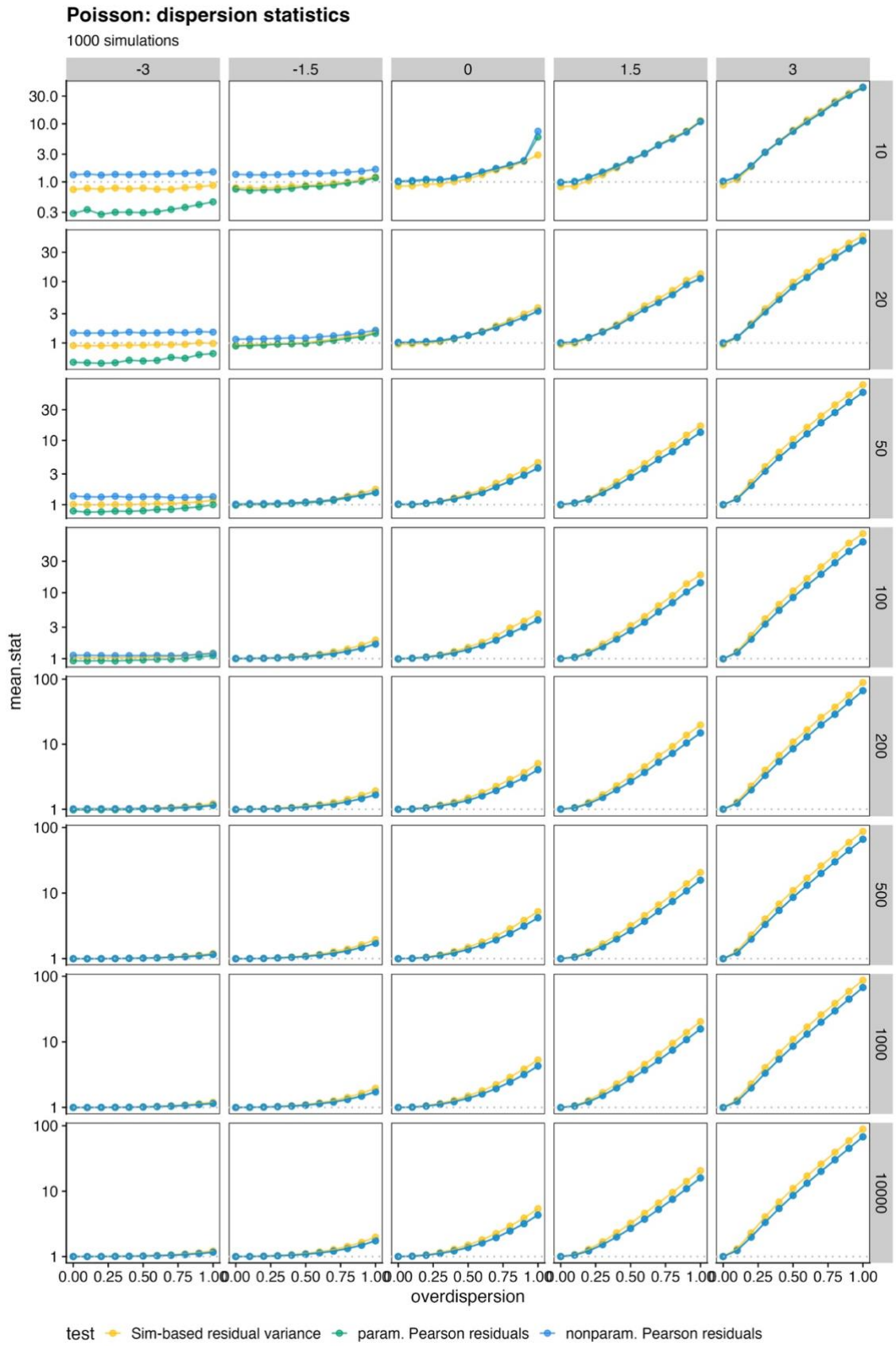
129 **S4. Dispersion statistics for GLMs**

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

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

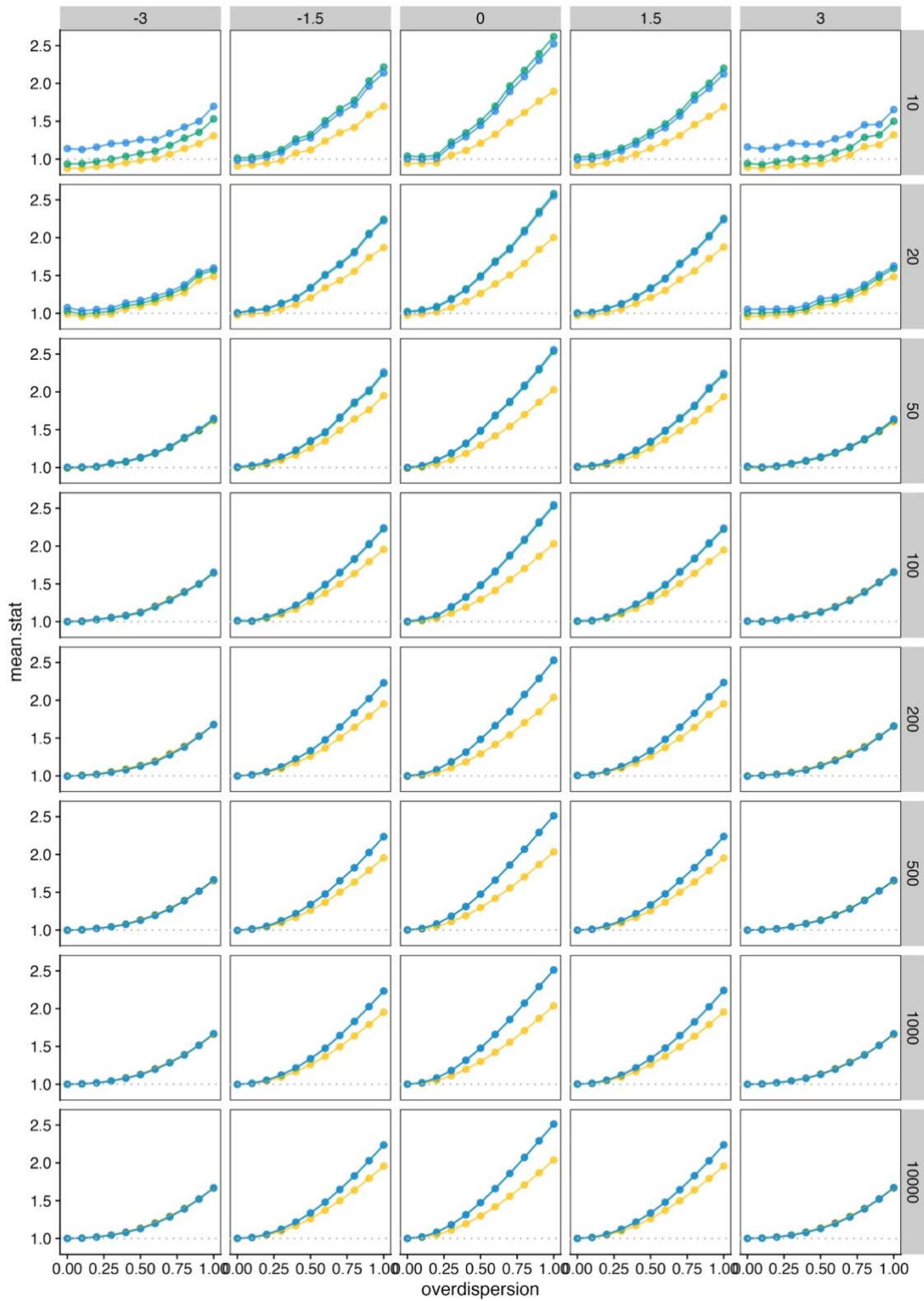


151

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

Binomial: dispersion statistics

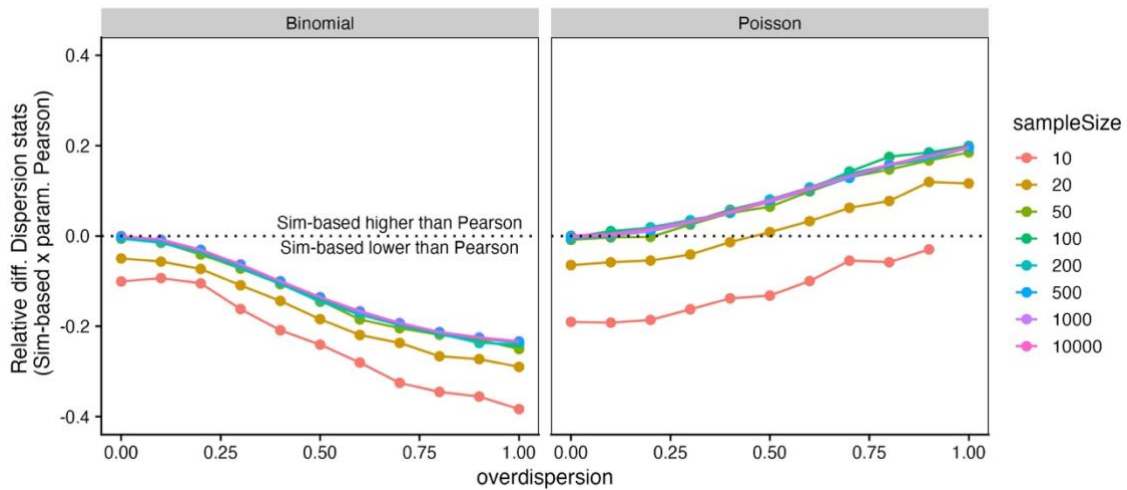
1000 sim; Ntrials=10



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

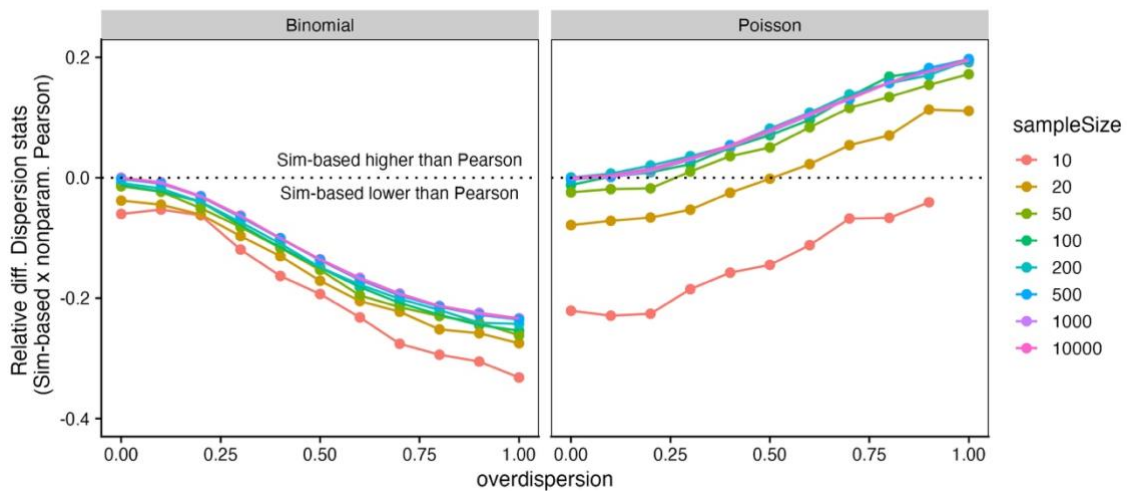
154

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



156

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

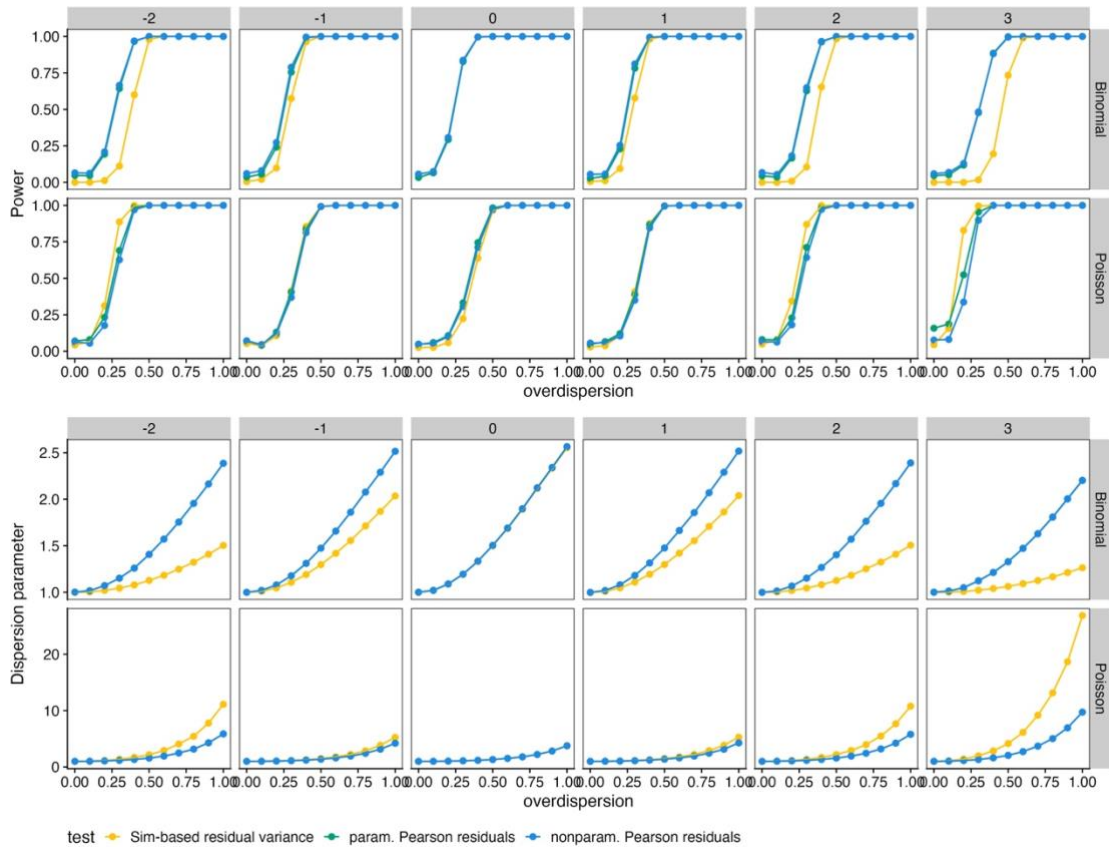


166

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

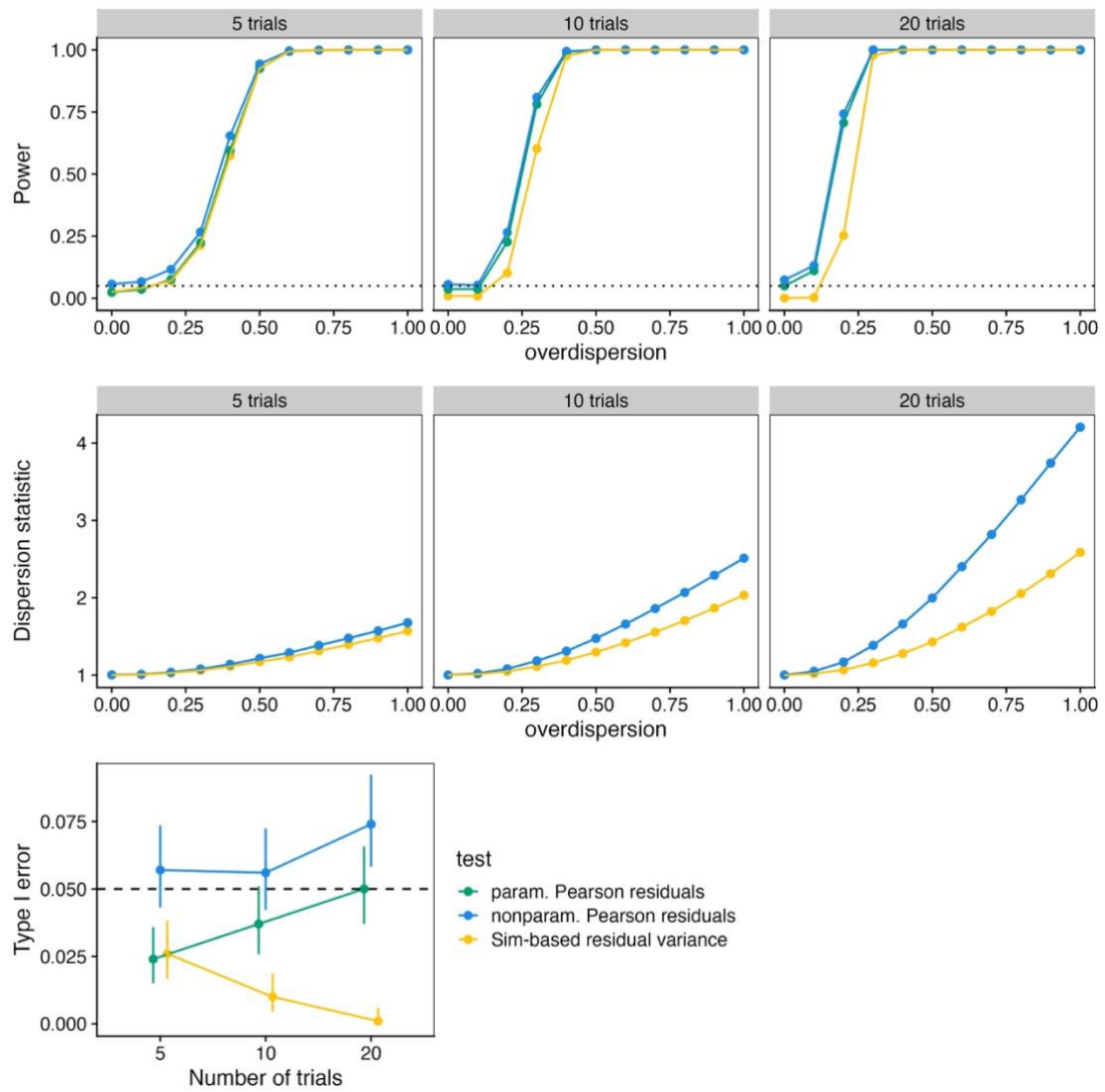
176 **S5: Expanding simulation parameters for GLMs**

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

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



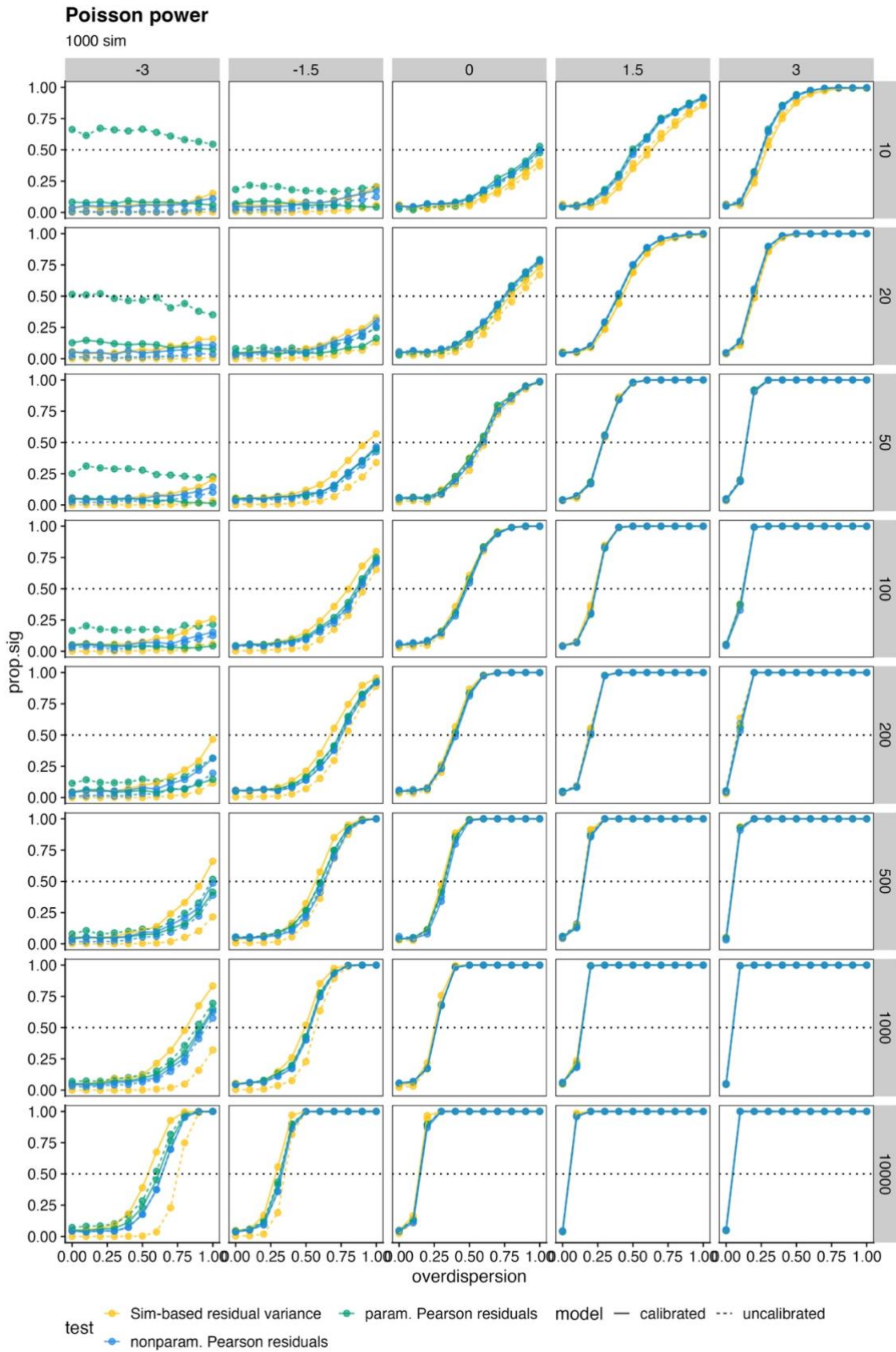
187

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

192 **S6. Power for the GLMs**

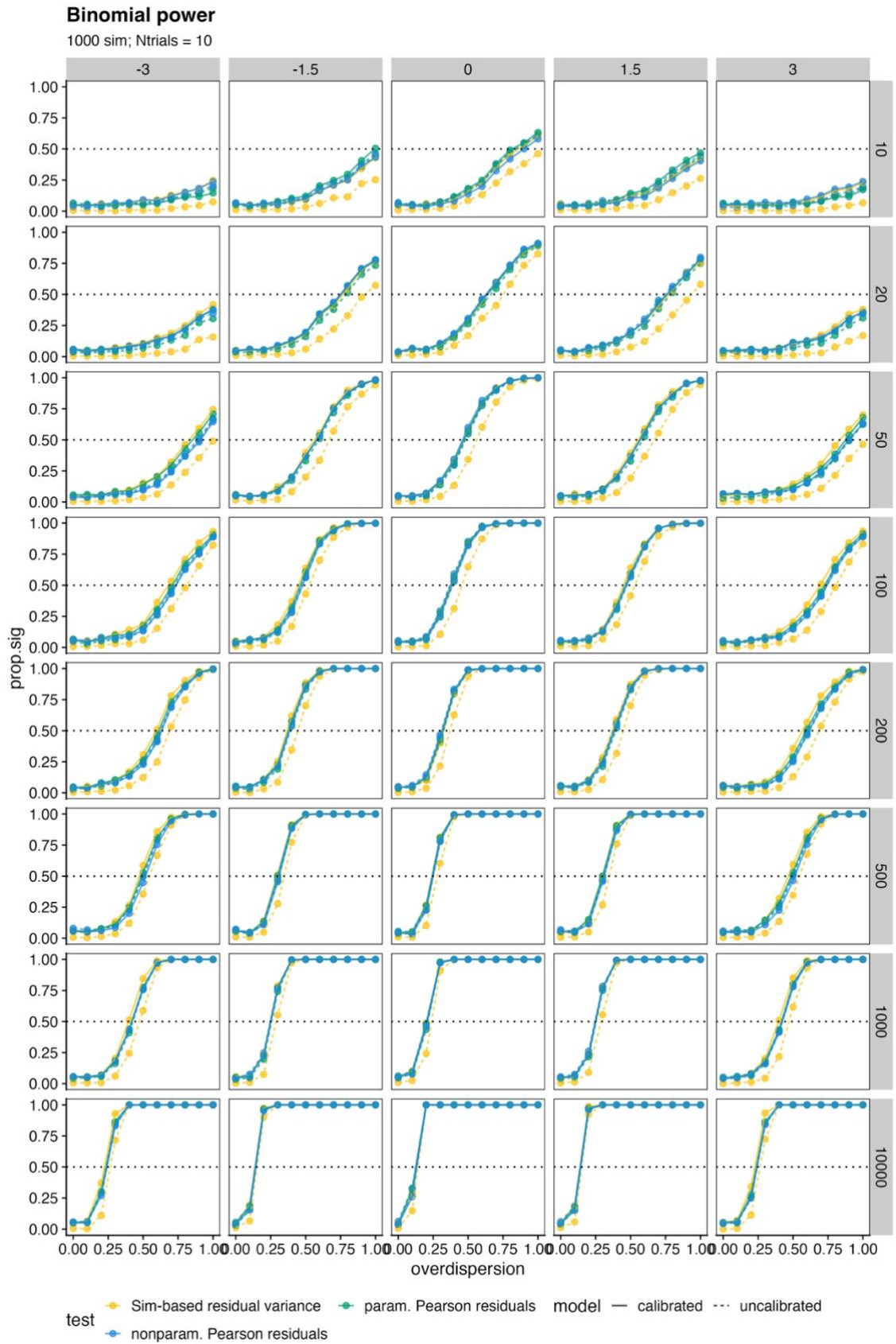
193 *Power calibration*

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



202

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



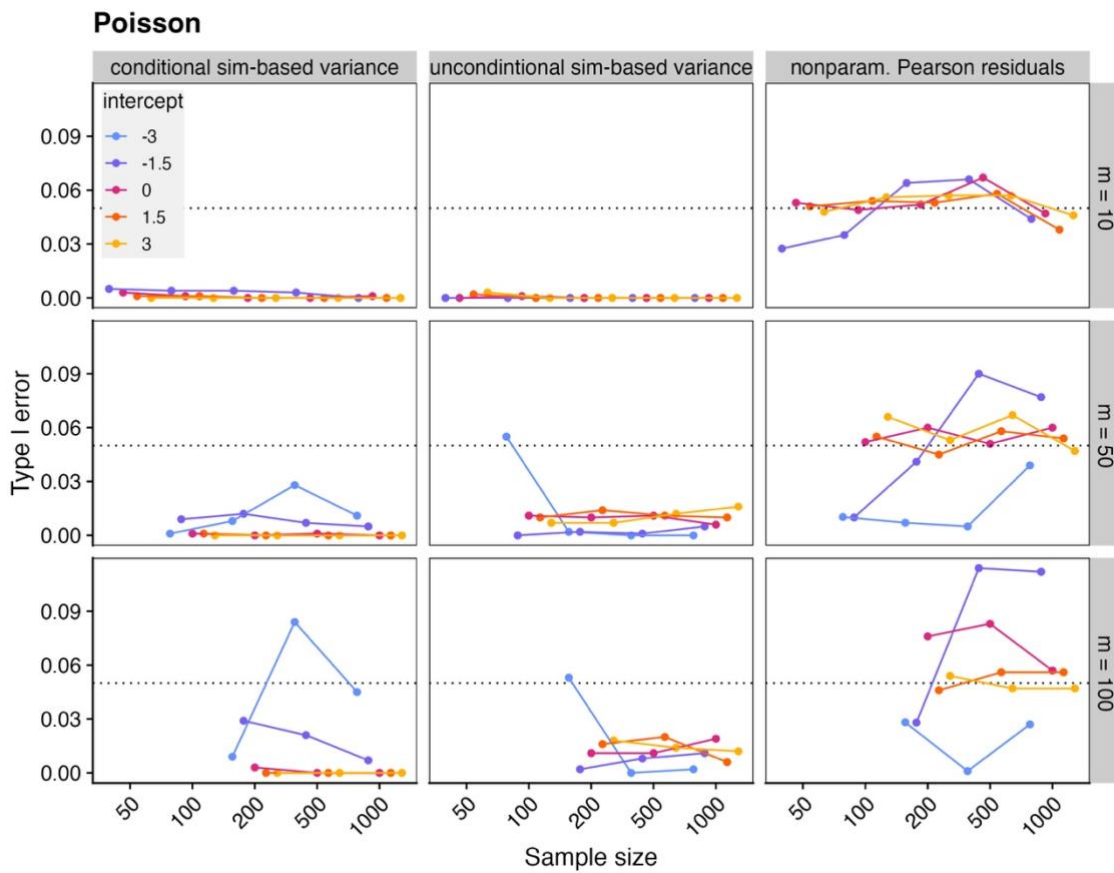
204

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

206 **S7. Additional GLMM results**

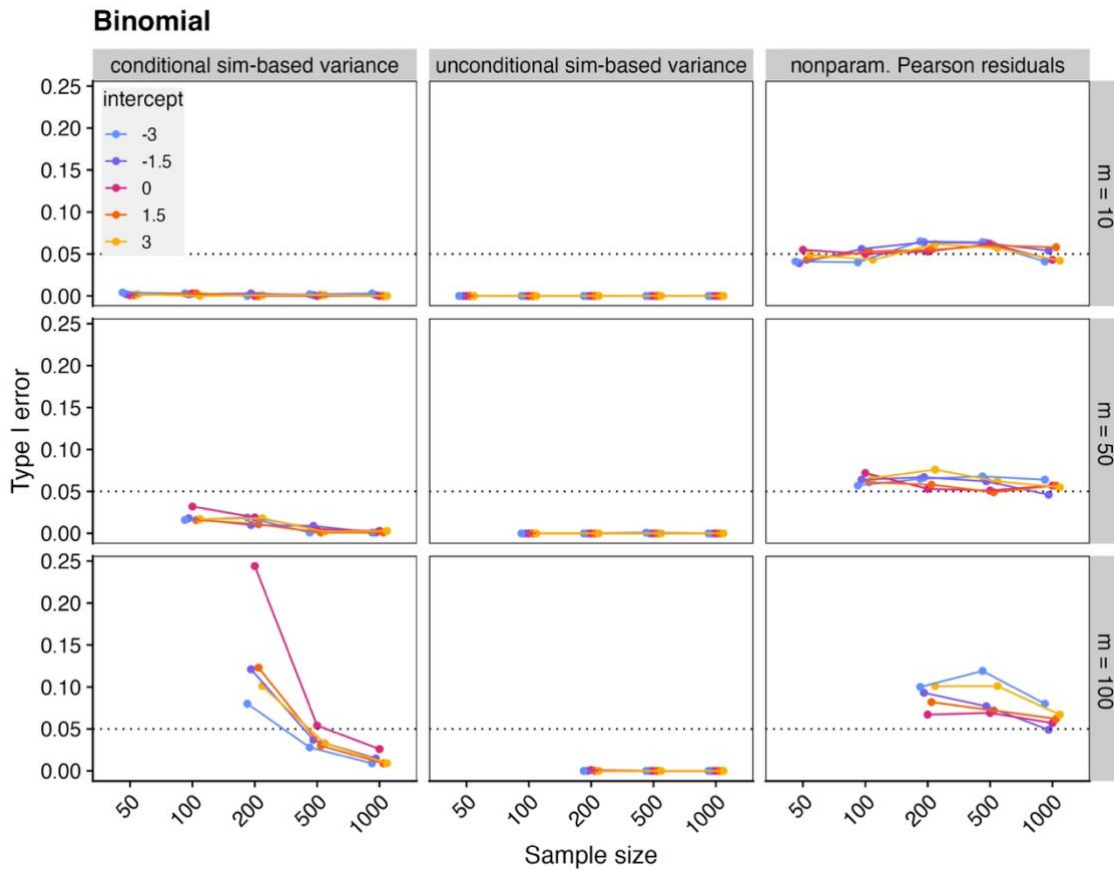
207 *Type I error rate of the alternative dispersion tests*

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

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

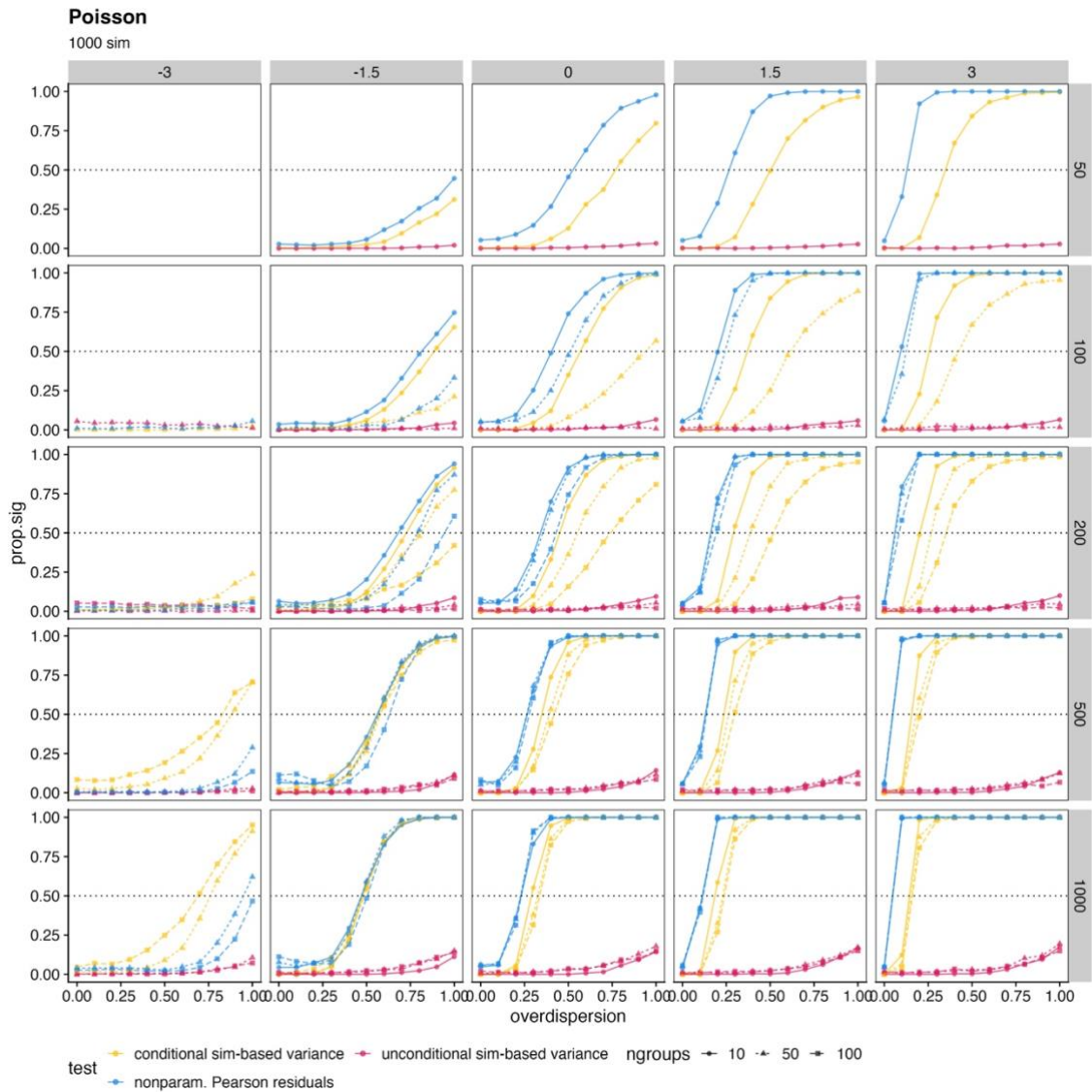


216

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

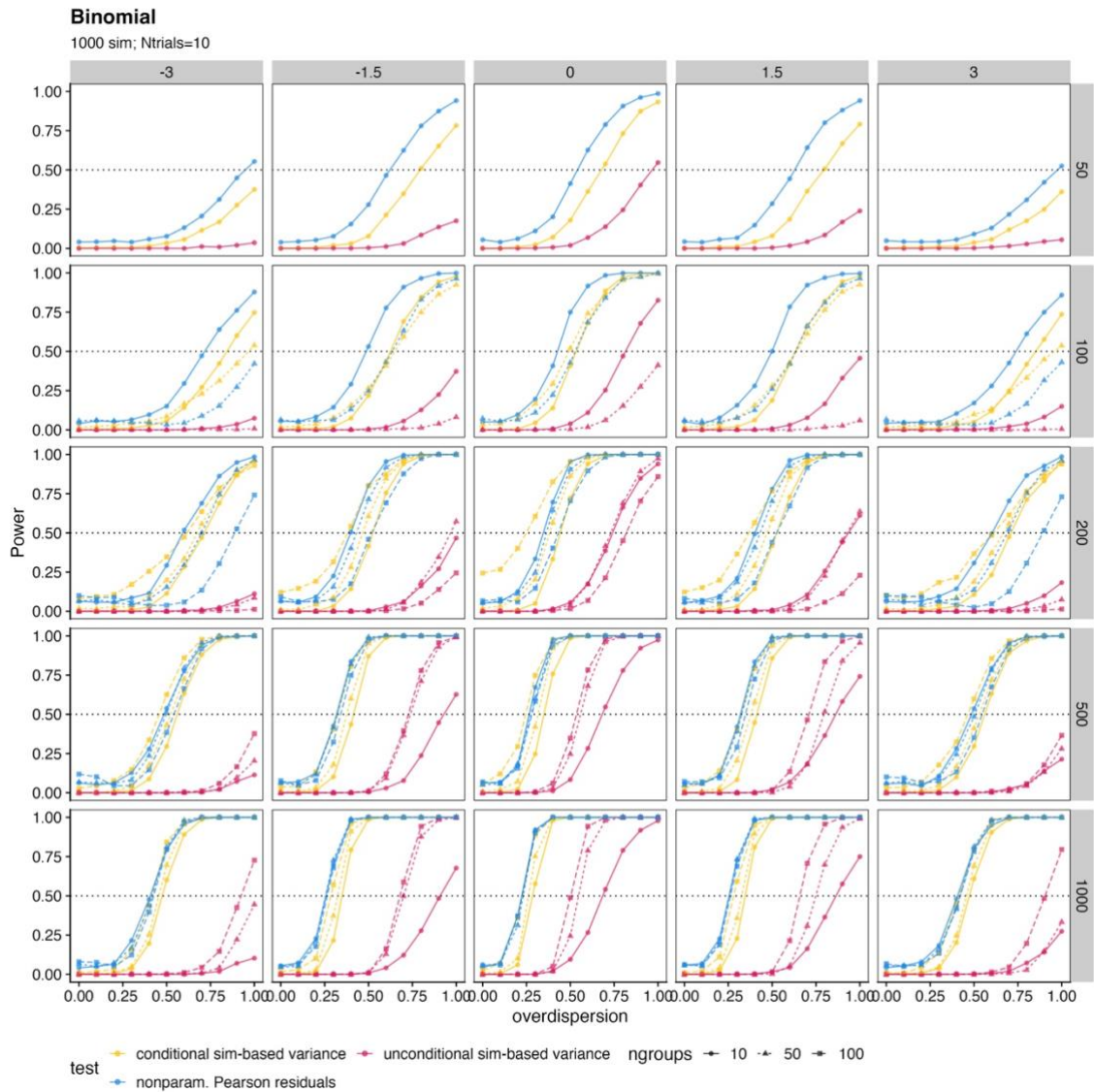
221 *Power of the alternative dispersion tests*

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



225

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

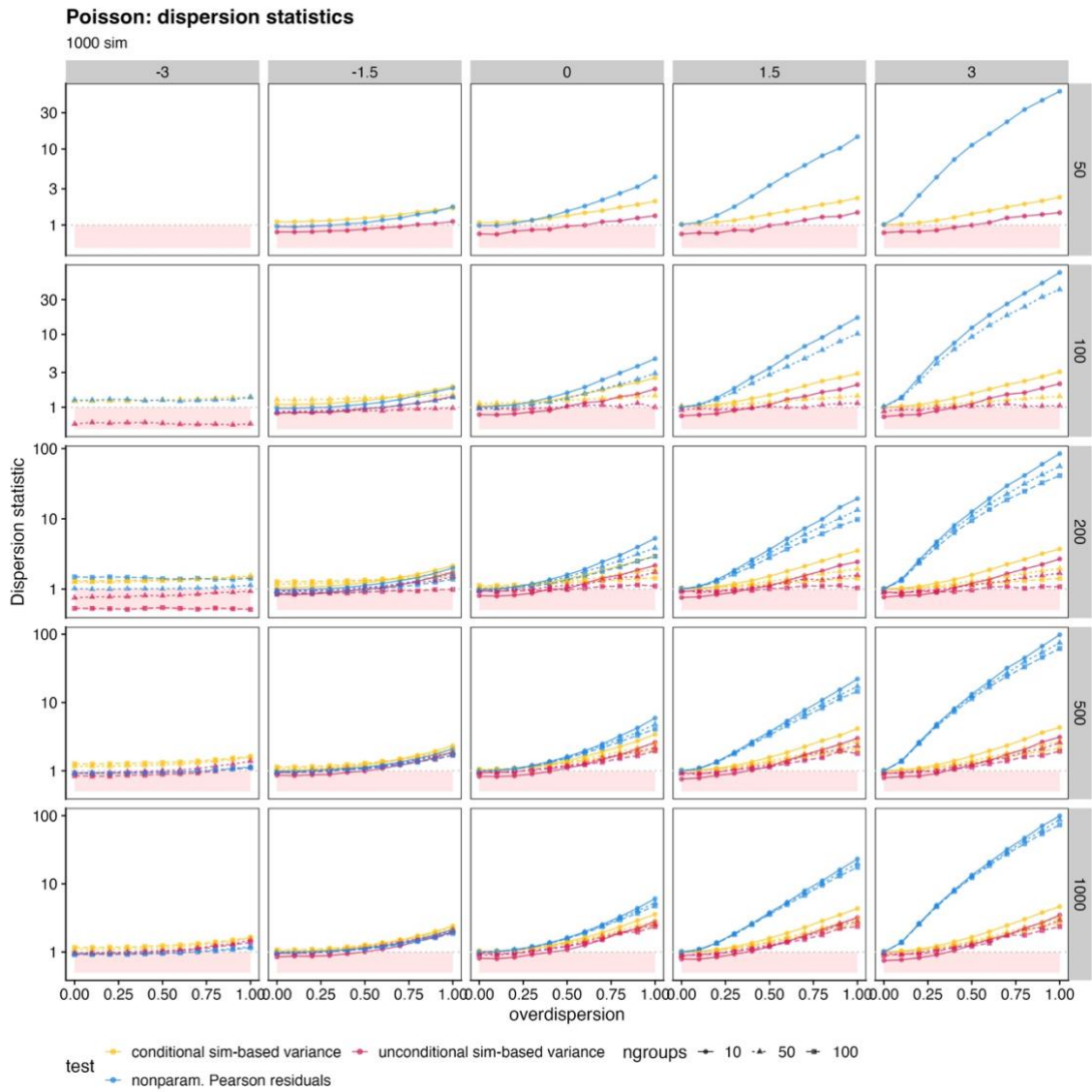


231

232 **Figure S7.4.** Power of the three alternative dispersion tests for binomial GLMMs, with
 233 different numbers of observations (rows), intercepts (columns), and number of groups
 234 for the random intercept (line types). 1000 simulations for each parameter set.

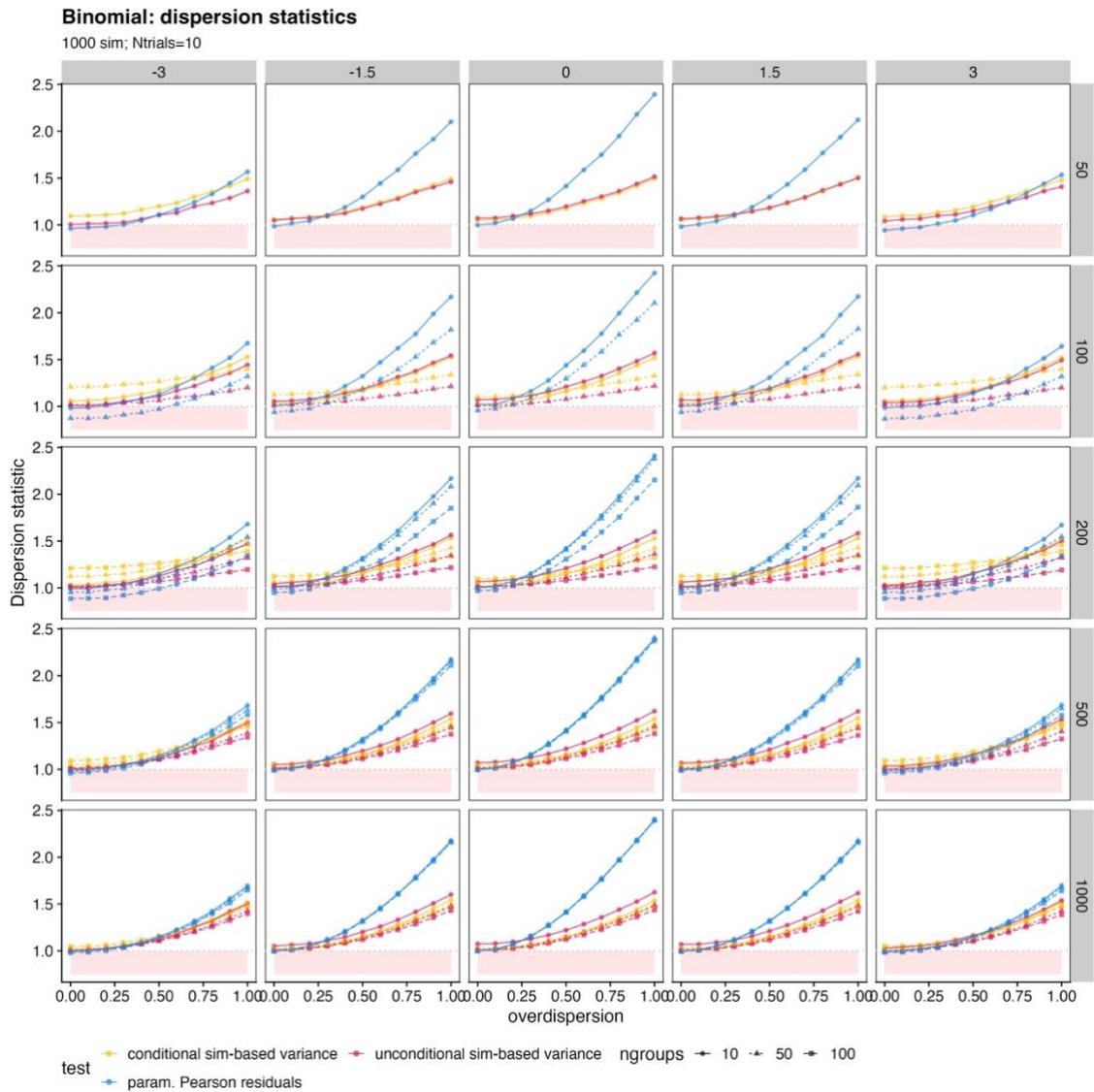
235 *Dispersion statistics of the alternative dispersion tests*

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



239

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



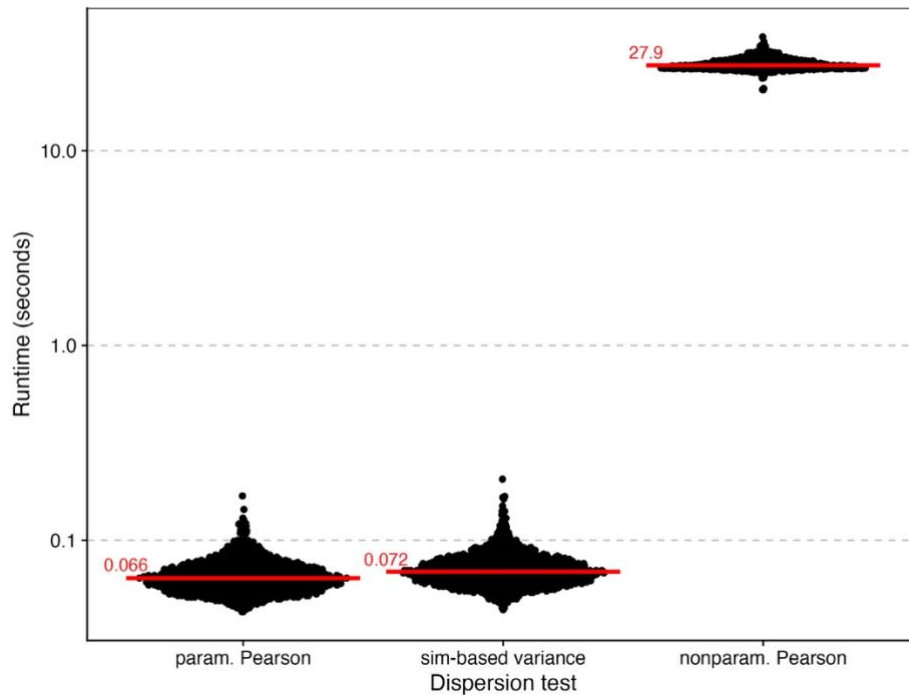
245

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

250 *Computational runtime for tests with GLMMs*

251 We computed the run time for the three tests used for GLMMs: the parametric
 252 Pearson test, the nonparametric Pearson test, and the simulation-based residual variance
 253 test with conditional simulations (Figure S7.7). We used 1,000 simulations of the
 254 Poisson GLMM as an example, with an overdispersion parameter of 0.4, an intercept of
 255 0, a sample size of 1,000, and 100 groups. There was almost no difference in

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

266 **S8: Alternative simulation-based residual variance test**

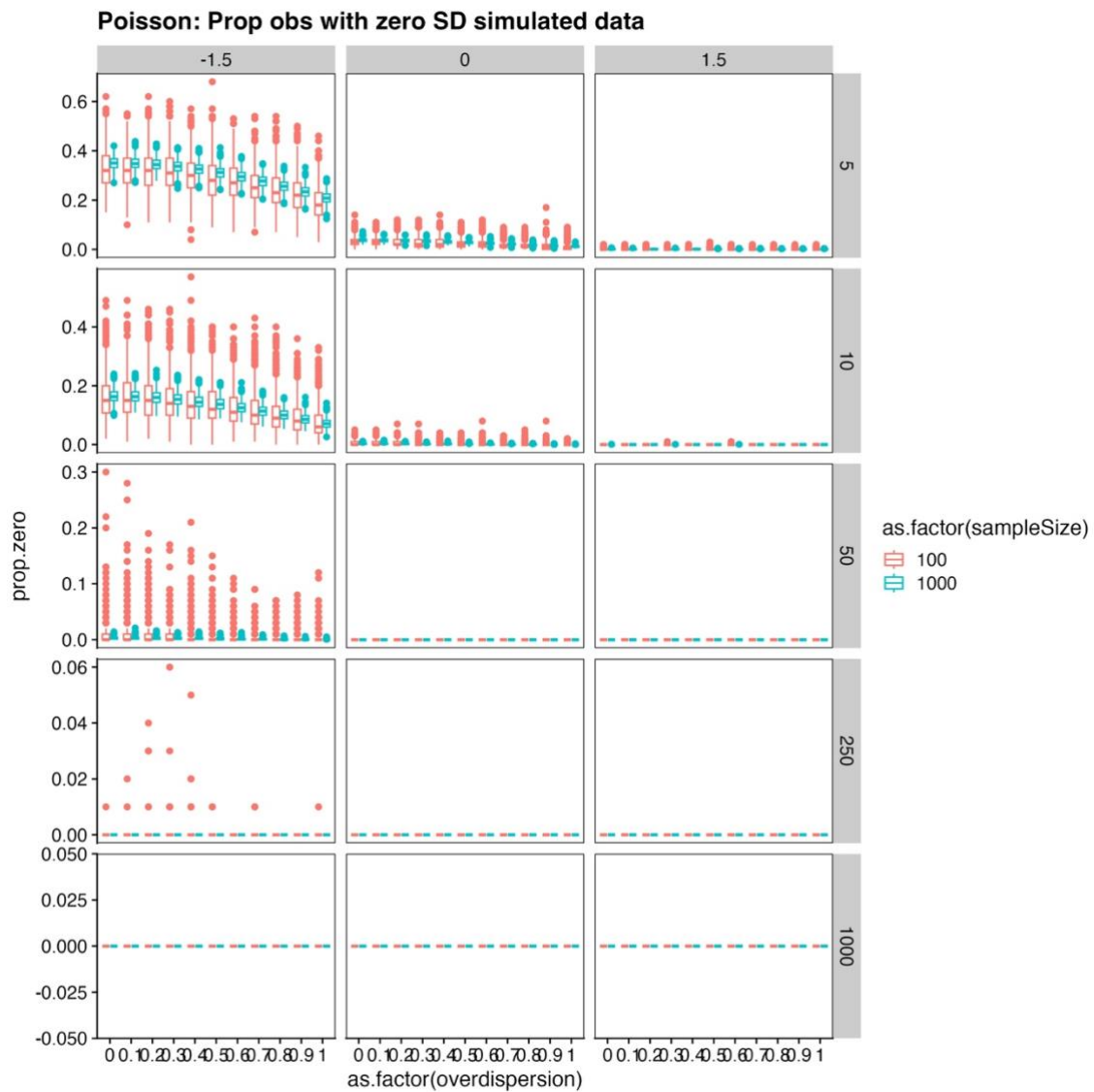
267 Another possibility for improving dispersion tests for GLMMs is to develop a
268 simulation-based approach that shows better type I, power, and a dispersion statistic that
269 could be interpreted similarly to the Pearson dispersion. To explore future possibilities,
270 we briefly considered an alternative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

$$275 \quad \text{Approx. Pearson observed residuals: } r_i = \frac{(y_i - \hat{\mu})}{\text{var}(y_{is})} \quad (\text{Equation S8.1})$$

$$276 \quad \text{Approx. Pearson simulated residuals: } r_{is} = \frac{(y_{is} - \hat{\mu})}{\text{var}(y_{is})} \quad (\text{Equation S8.2})$$

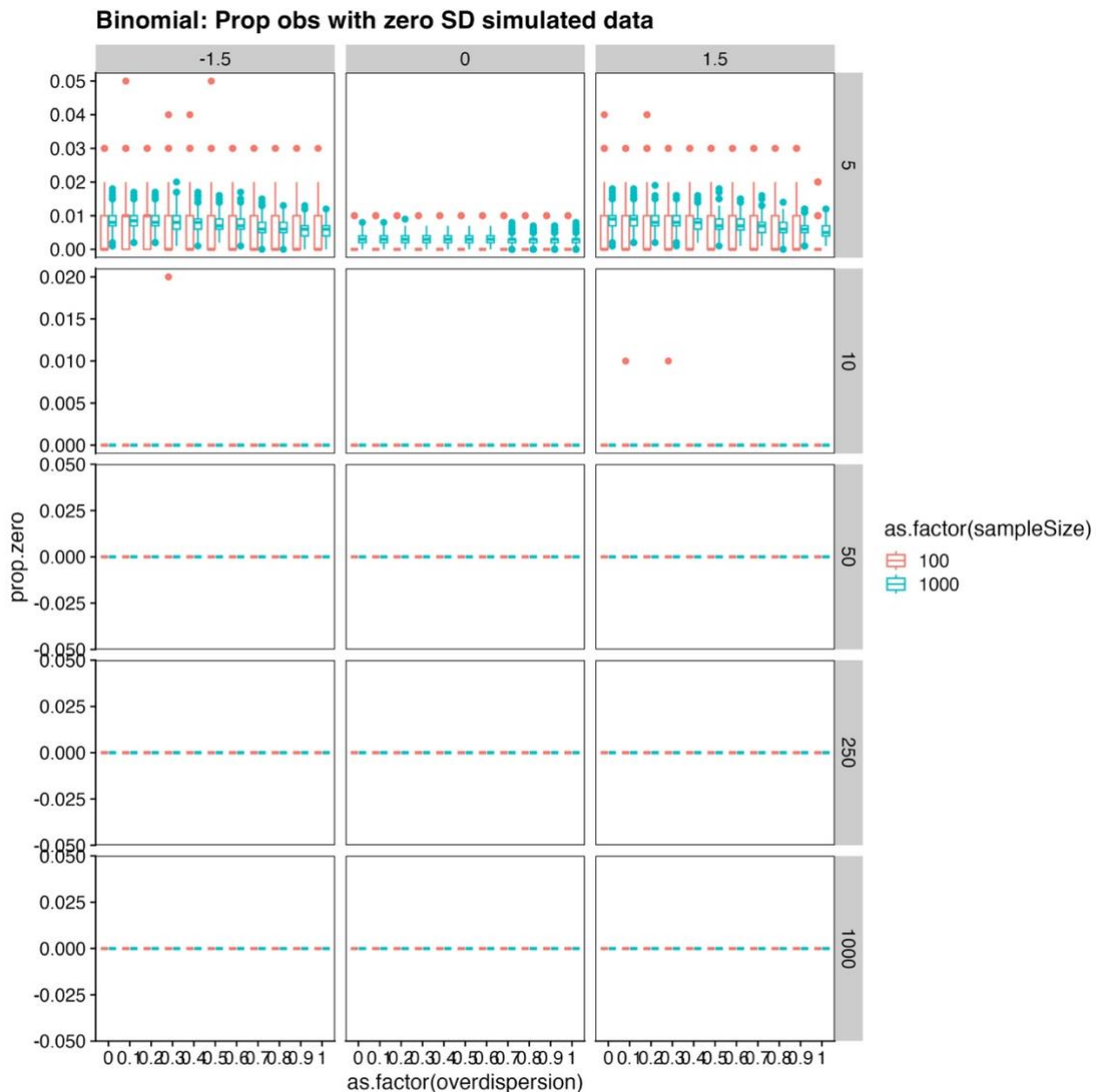
277 One obstacle with calculating the denominator of the approximate Pearson
278 residuals for each observation is that the variance depends on the number of simulations
279 and the model parameters, such as the intercept or the number of trials in the binomial
280 GLM/GLMMs. If there are too few simulations or the intercept is very small, the
281 chance of resulting in zero variance (all simulated values are the same) is higher for data
282 points with small variance. To overcome this, we first evaluated the minimum number
283 of simulations for different intercepts and sample sizes, in which all observations have
284 estimated variances that are different from zero. For all combinations of parameters, we
285 found that 1,000 simulations were sufficient to ensure that all variances in the simulated
286 observations were positive (Figures S8.1 and S8.2). However, 250 simulations (the
287 default parameter of the DHARMA package) also presented reasonable results, with the
288 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

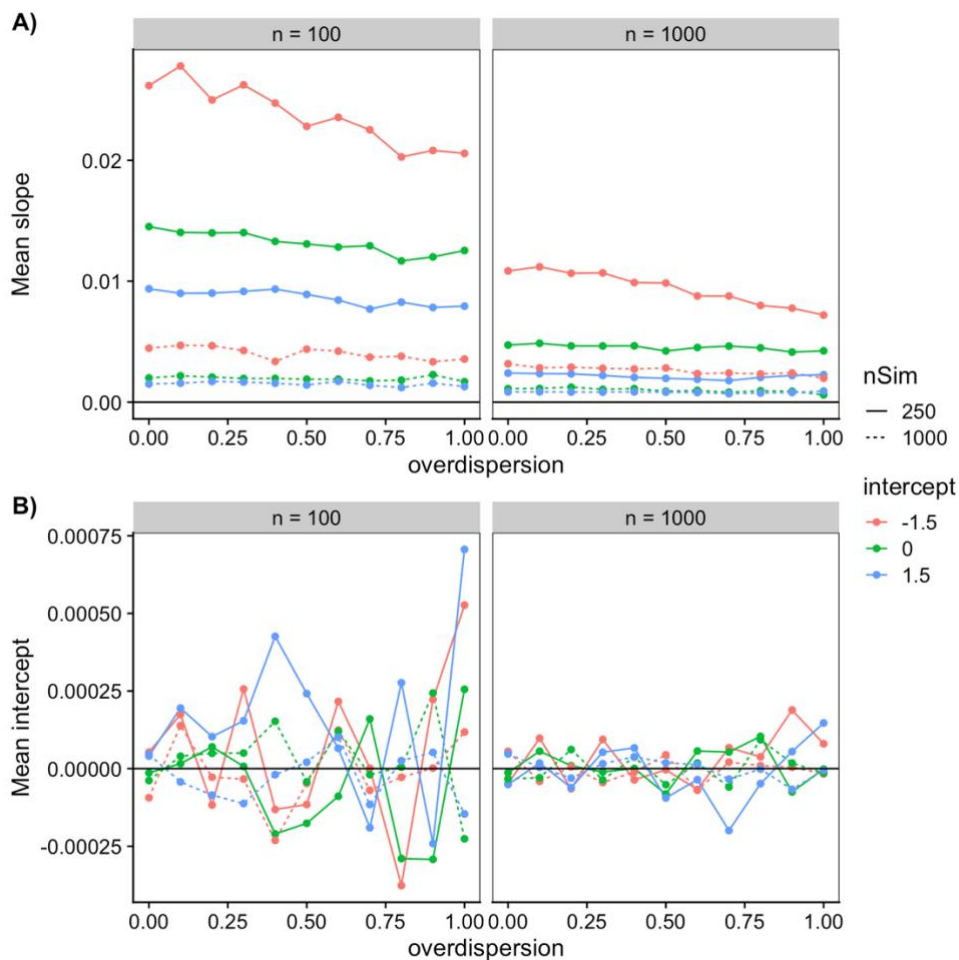
298 **Figure S8.1.** Poisson GLM: Proportion of observations with simulated zero variance in
 299 the dataset for different combinations of intercept (columns), number of simulations
 300 (rows) and sample sizes (colours).



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

306 First, we compared the approximate Pearson residuals for GLMs with the
 307 Pearson residuals by regressing the difference between them as the response variable
 308 and the Pearson residuals as the predictor for the Poisson GLMs (Figure S8.3). The
 309 intercepts for all simulation sets were nearly zero. The slope of the regression was
 310 positive and very small for the larger number of simulations and intercepts. It means

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

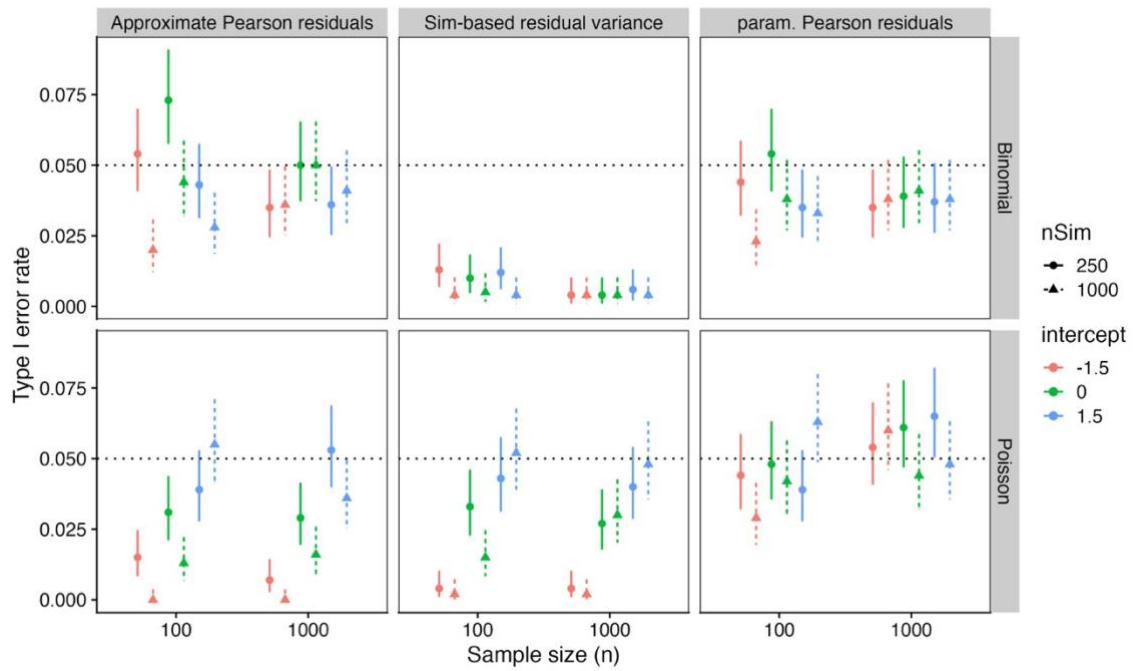


313

314 **Figure S8.3.** Mean slope (A) and intercept (B) of the regression of the difference
 315 between the Approximate Pearson residuals and Pearson residuals as response variable
 316 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

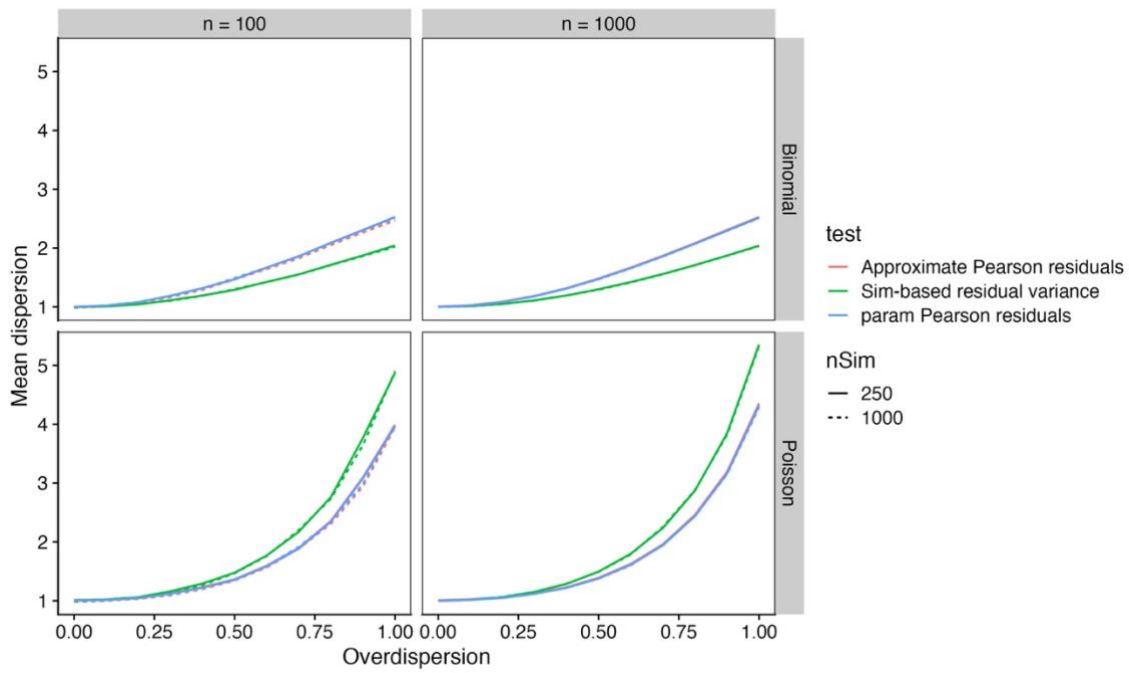


323

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

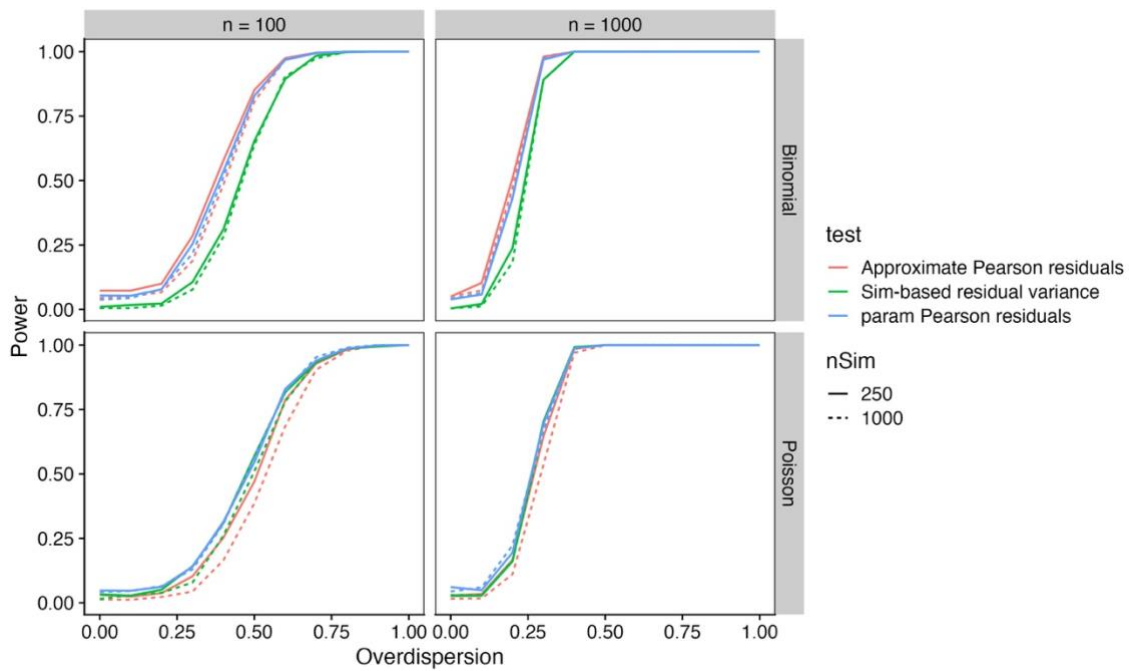
327

328 The dispersion statistics for the alternative simulation-based residual variance
329 test didn't change depending on the number of simulations and were very similar to the
330 parametric Pearson dispersion statistics for both GLMs (Figure S8.5). Power was very
331 similar among the tests for the Poisson GLM (Figure S8.6). For binomial GLMs, the
332 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.

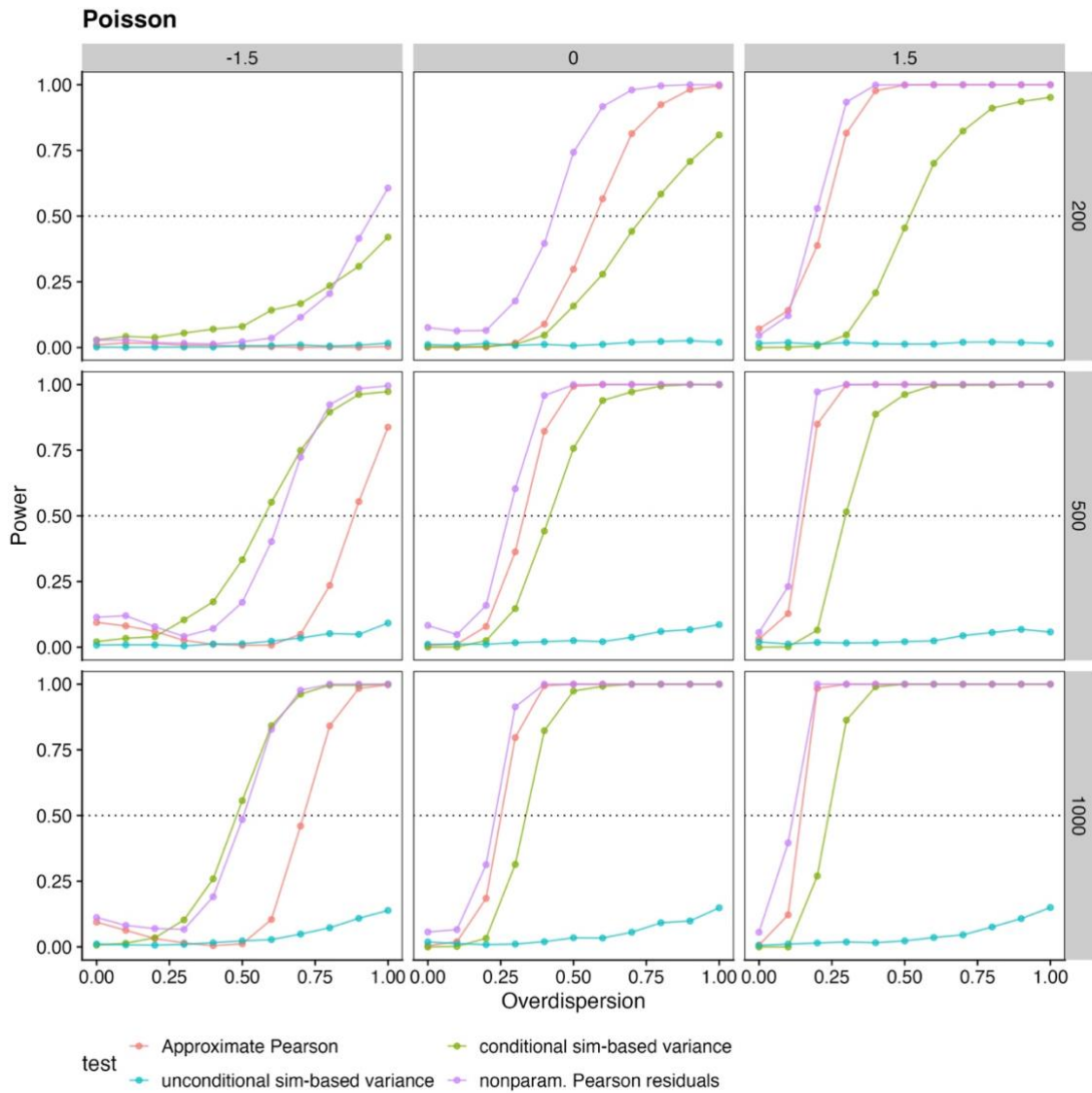


335

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

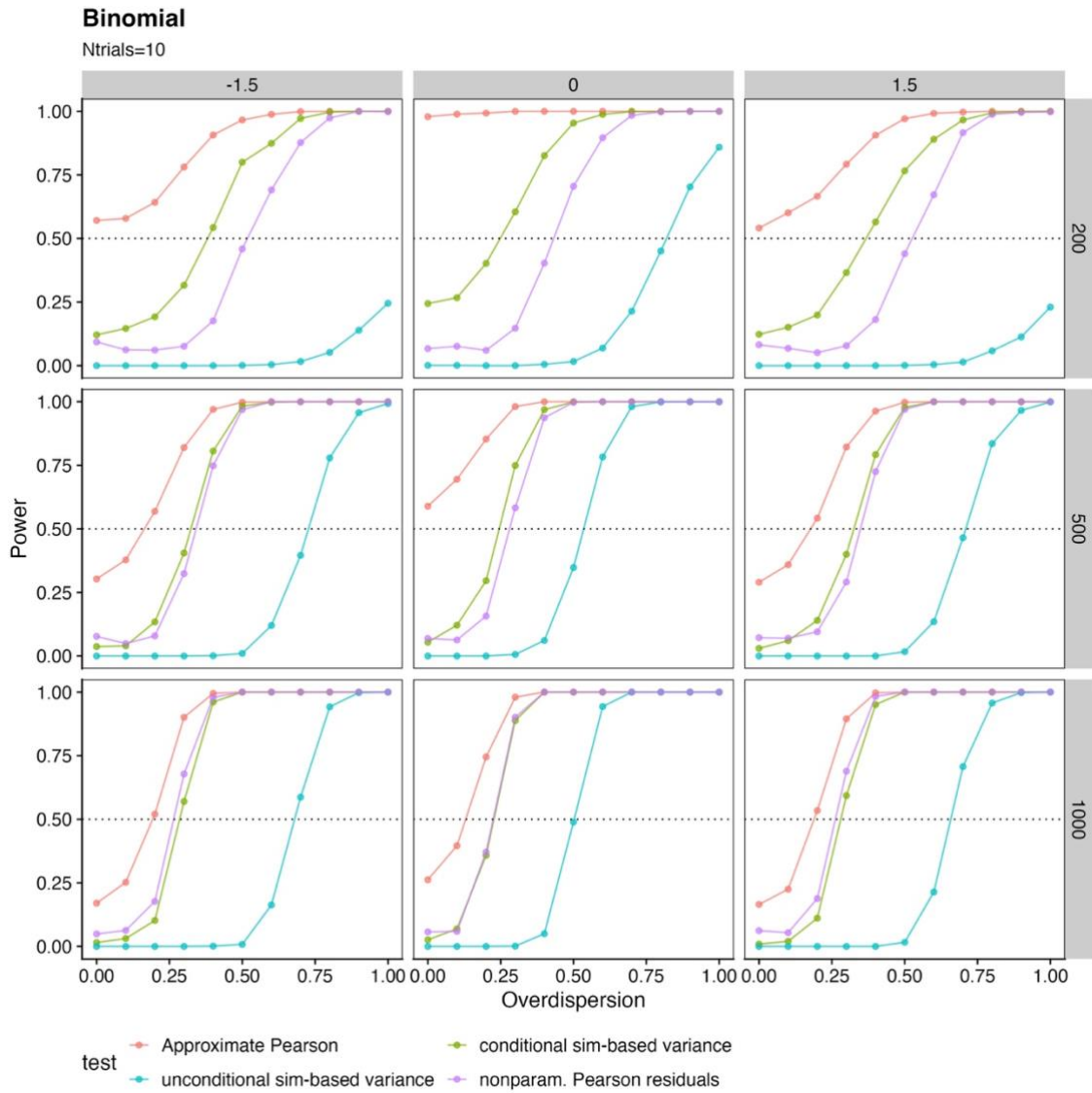
337 For the GLMM simulations, we fixed the number of groups at 100 and the
 338 number of simulations at 250 to compare with the cases where the Pearson Chi-squared
 339 test fails. We compared sample sizes of 200, 500, and 1000 observations and intercepts
 340 of -1.5, 0, and 1.5. We excluded simulations 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.



350

351 **Fig S8.8.** Power for binomial GLMMs for the alternative simulation-based test using an
 352 approximation for Pearson residuals compared with the other tests assessed in the study.
 353 1000 simulations for each parameter set: intercept (panel columns) and sample size
 354 (panel rows). The fixed parameters are slope = 1, number of groups = 100, random
 355 effects variance = 1, number of trials = 10.

356 **S9. Parametric Pearson test with approximated residual degrees of**
357 **freedom for GLMMs**

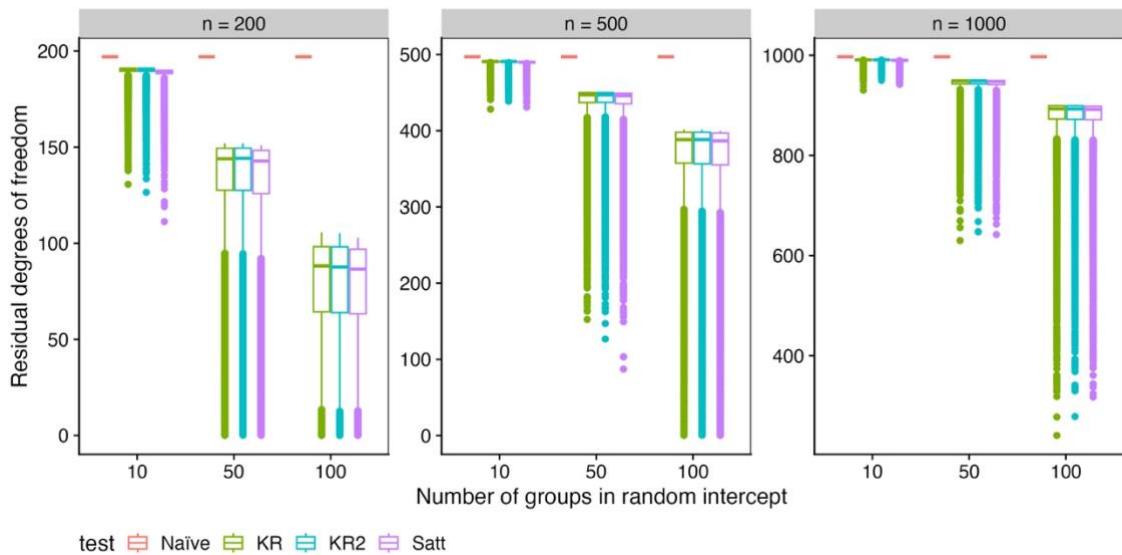
358 Degrees of freedom (*df*) are not always known for GLMMs with complex
359 hierarchical structures and limit the use of the parametric Pearson test because it
360 depends on it for evaluating overdispersion with the Chi-squared distribution.
361 Moreover, our results show that using the naïve *df* is problematic for testing dispersion
362 when you have a large number of groups in the random intercept. The two most
363 suggested methods to approximate *df* of mixed-effect models, the Satterthwaite (1946)
364 and the Kenward-Roger (Kenward & Roger 2009), were developed for LMMs to
365 account for the effect of the covariance structure on *df* and standard errors. Stroup et al.
366 (2013) suggested that the adjustment is also accurate for GLMMs. However, none of the
367 most used R packages use any correction for the degrees of freedom for GLMMs. The
368 few R packages that provide those approximations, e.g. *lmerTest* (Kuznetsova et al.,
369 2017; Kuznetsova et al., 2020) that relies on *pbkrtest* (Halekoh & Højsgaard 2014), are
370 only implemented for LMMs.

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

- 375 - The Kenward-Roger (KR) bias-corrected variance-covariance matrix for the
376 fixed effect parameters and degrees of freedom from Kenward & Roger (1997).
- 377 - The improved correction of the Kenward-Roger (KR2) returns an improved
378 correction given in Kenward & Roger (2009).
- 379 - The Satterthwaite correction (Sat) from Satterthwaite (1946).

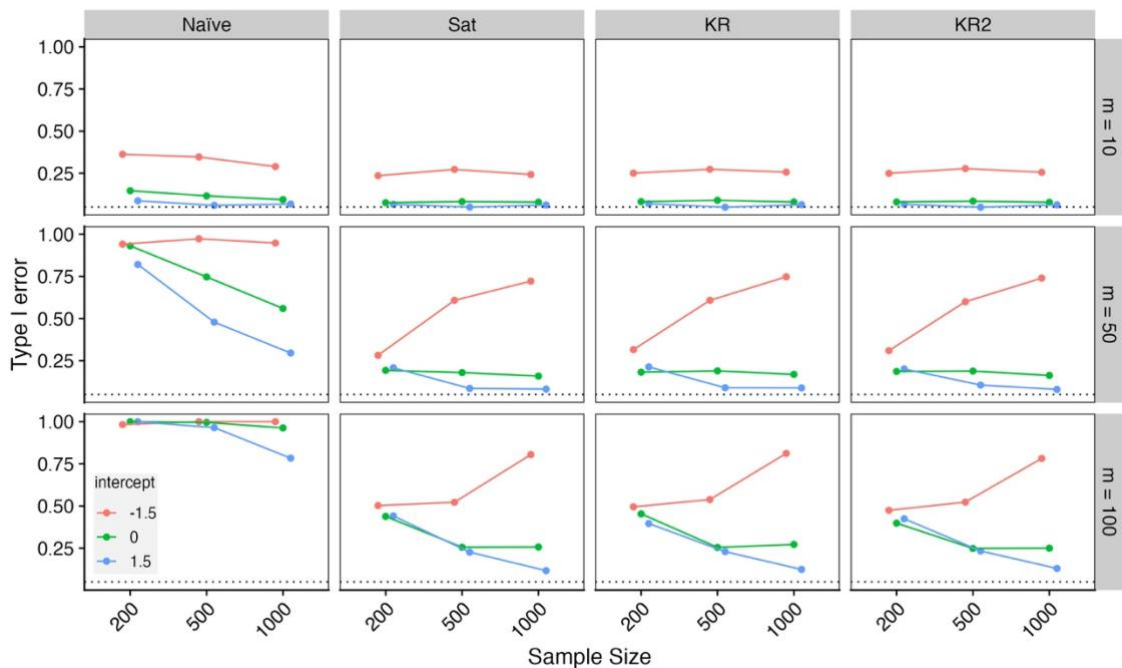
380 Our test results show that all three correction methods presented very similar
381 residual df for all simulation settings (Figure S9.1), which resulted also in very similar
382 test results (e.g., Figure S9.2 for type I error). Given the high similarity among tests for
383 the different residual 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.



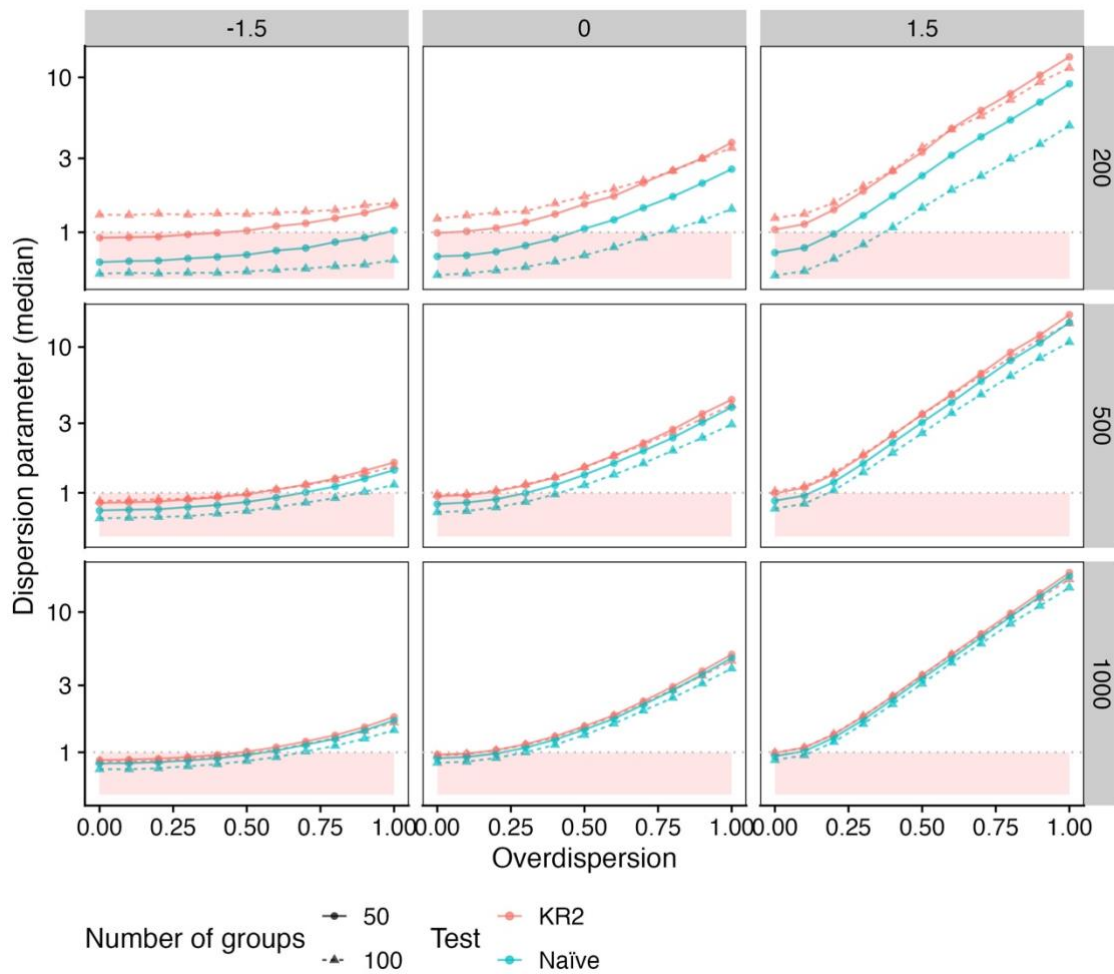
396

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



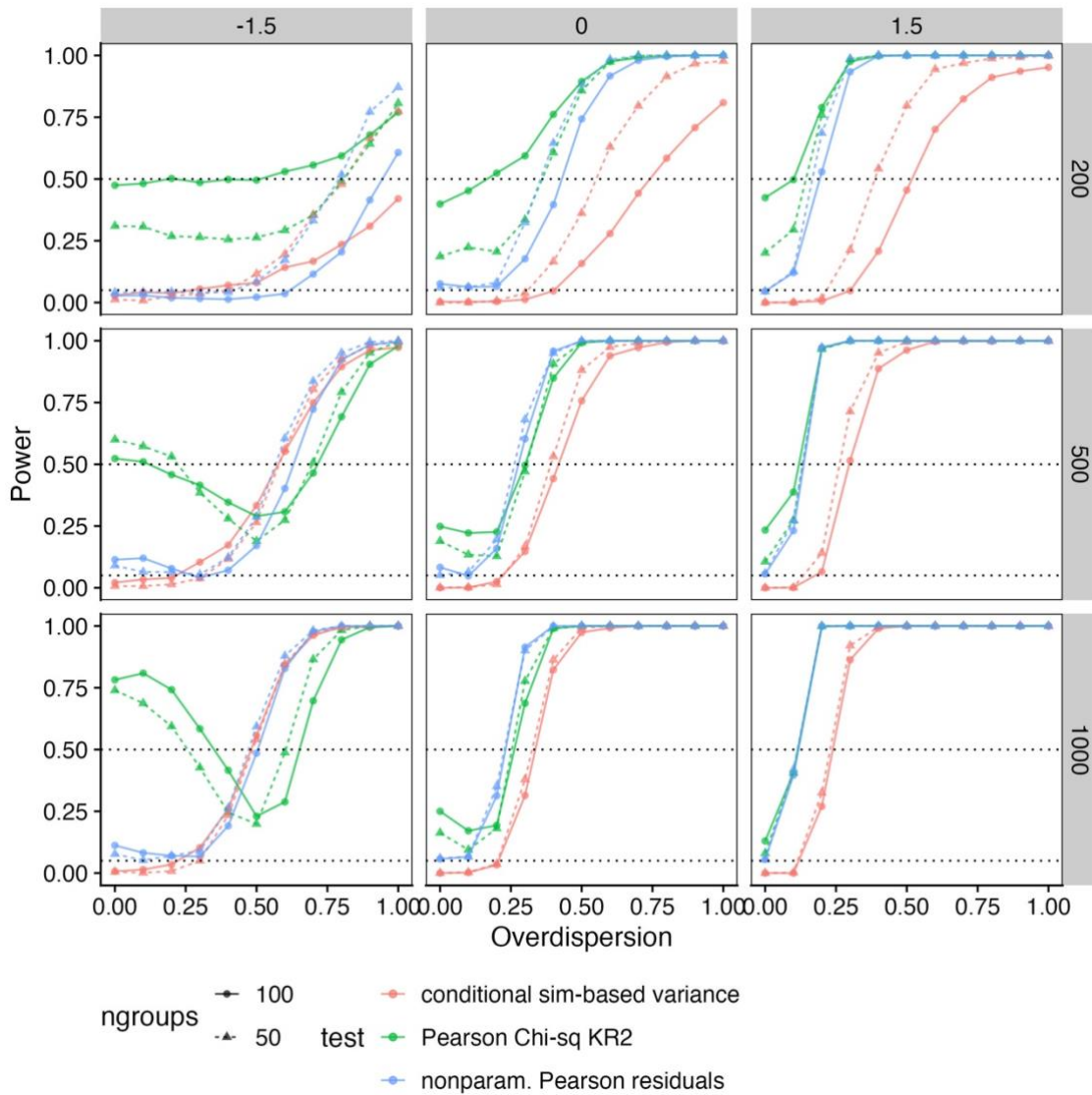
402

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



409

410 **Figure S9.3.** Dispersion parameters for the parametric Pearson test for Poisson GLMMs
 411 performed with different corrections for the residual degrees of freedom (colours),
 412 number of groups in the random intercept (linetype and shape), sample size (panel
 413 rows), and intercept (panel columns). Please refer to the main text above to relate to
 414 each applied correction. To improve clarity, we omitted the other corrections because
 415 they are too similar to each other. 1000 simulations for each parameter setting, slope =
 416 1, random intercept variance = 1.



417

418 **Figure S9.4.** Power of dispersion tests for Poisson GLMMs (colours) performed with
 419 different numbers of groups in the random intercept (linetype and shape), sample size
 420 (panel rows), and intercept (panel columns). Please refer to the main text above to relate
 421 to the applied correction for residual degrees of freedom. To improve clarity, we omitted
 422 other corrections for residual degrees of freedom because they are too similar to each
 423 other. 1000 simulations for each parameter setting, slope = 1, random intercept variance
 424 = 1.

425 **S10: Case studies detailed information**

426 *Case study 1: Redstart breeding pairs*

427 Summary results of the models applied to the Common redstart breeding pairs in
 428 Switzerland. Table S10.1 shows coefficients of models applied using a Poisson GLM
 429 (*glm* function in base R), a spatial explicit Poisson GLM (exponential spatial
 430 autocorrelation, *glmmTMB* R package) and a negative binomial GLM (*glm.nb* in MASS
 431 Rpackage). Table S10.2 shows the three dispersion tests (*testDispersion* in DHARMA)
 432 and Table S10.3 the Moran’s I residual spatial autocorrelation tests
 433 (*testSpatialAutocorrelation* in DHARMA) applied to all models.

434 **Table S10.1.** Fixed effect coefficients at the link scale (log) of the models fitted to the
 435 redstarts dataset. The effect of interest (forest cover) is highlighted in bold.

Model	Term	Estimate	SE	Statistic	P-value
Poisson	Intercept	0.137	0.105	1.314	0.189
	elevation	0.147	0.093	1.578	0.115
	elevation^2	-0.238	0.075	-3.156	0.002
	forests	-0.007	0.003	-2.213	0.027
negative binomial	Intercept	0.145	0.164	0.882	0.378
	elevation	0.106	0.138	0.77	0.441
	elevation^2	-0.231	0.101	-2.284	0.022
	forests	-0.007	0.005	-1.61	0.107
spatial Poisson	Intercept	-0.479	0.235	-2.037	0.042
	elevation	0.196	0.174	1.121	0.262
	elevation^2	-0.239	0.109	-2.203	0.028
	forests	-0.006	0.005	-1.429	0.153

436

437

438 **Table S10.2.** Residuals dispersion tests of the models fitted to the redstarts dataset. The
 439 parametric Pearson test for the spatial Poisson GLM is just for overdispersion, and the
 440 very small dispersion coefficient happens because we used a GLMM structure to model
 441 it.

Model	Test	Dispersion	P-value
Poisson	sim-based res. variance	2.47	0
	parametric Pearson	2.38	0
	nonparametric Pearson	2.39	0
negative binomial	sim-based res. variance	1.03	0.808
	parametric Pearson	1.02	0.775
	nonparametric Pearson	1.04	0.472
spatial Poisson	sim-based res. variance	0.81	0.936
	parametric Pearson	0.39	1
	nonparametric Pearson	0.99	0.984

442

443 **Table S10.3** Moran's I residual spatial autocorrelation of the models fitted to the
 444 redstart dataset.

Model	Moran's I	Expected	SD	P-value
Poisson	0.0258	-0.0029	0.0043	0
negative binomial	0.0078	-0.0029	0.0043	0.0128
spatial Poisson	-0.0003	-0.0029	0.0043	0.5407

445

446 *Case study 2: Wild and zoo-housed orangutan behavior*

447

448 **References**

- 449 Hartig, F. (2024). *DHARMA: Residual Diagnostics for Hierarchical (Multi-Level /*
450 *Mixed) Regression Models* (Version 0.4.7) [Computer software].
451 <https://CRAN.R-project.org/package=DHARMA>
- 452 Jahn, N. (2023). *europemc: R interface to the europe PubMed central restful web*
453 *service* (Version 0.4.3) [Computer software]. [https://CRAN.R-](https://CRAN.R-project.org/package=europemc)
454 [project.org/package=europemc](https://CRAN.R-project.org/package=europemc)
- 455 Laumer, I. B., Kansal, S., Van Cauwenberghe, A., Rahmaeti, T., Setia, T. M., Mundry,
456 R., Haun, D., & Schuppli, C. (2025). Wild and zoo-housed orangutans differ in
457 how they explore objects. *Scientific Reports*, *15*(1), 14853.
458 <https://doi.org/10.1038/s41598-025-97926-z>
- 459 Lüdecke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021).
460 performance: An R package for assessment, comparison and testing of statistical
461 models. *Journal of Open Source Software*, *6*(60), 3139.
462 <https://doi.org/10.21105/joss.03139>
- 463

1 **Supporting information for:**

2 **Dispersion tests in generalized linear mixed-effects models: a**

3 **methods comparison and practical guide for ecologists**

4

5 **S1. Trend analysis and current ecological literature practices on**

6 **dispersal issues**

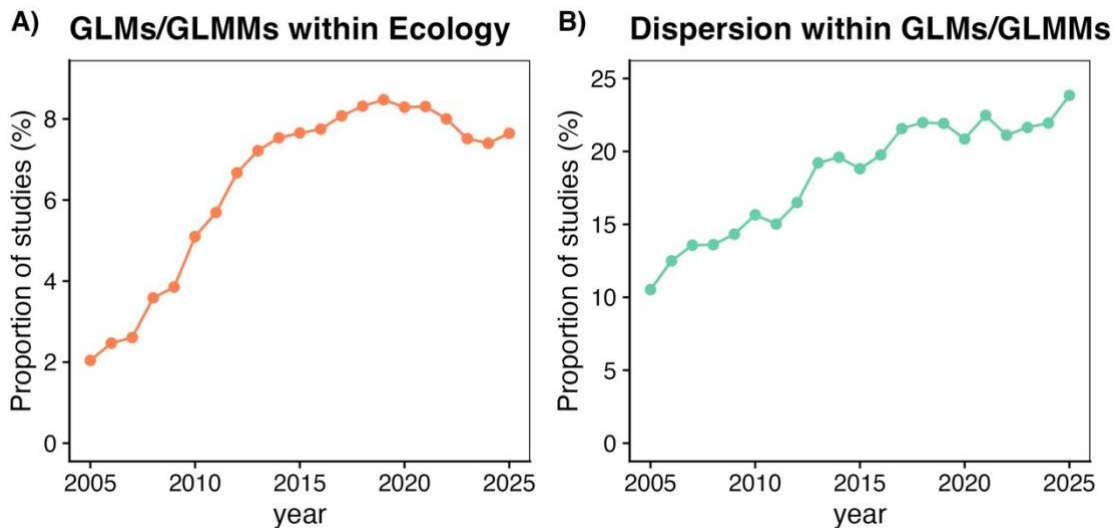
7 To understand the extent of ecological studies relying on GLMs/GLMMs for count and
8 discrete proportion data and those that address dispersion issues, we conducted a text analysis of
9 the ecological literature over the past 20 years. We used the R package ‘europepmc’ (v0.4.3,
10 Jahn, 2023) to search for articles in the PubMed and Medline NLM databases from 2005 to
11 2025. We used combinations of words (Table S1.1) to retrieved the annual records for: (1)
12 the percentage of ecological papers using GLMs/GLMMs for count and discrete
13 proportion data (Figure S1.1a), (2) the percentage of those papers that mention
14 dispersion terms in general (Figure S1.1b), (3) the percentage of ecological papers using
15 GLMs/GLMMs for count data mentioning dispersion terms (Figure BOX 1, main text),
16 and (4) the percentage of ecological papers using GLMs/GLMMs for discrete
17 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.



27

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

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

44 Of the 40 papers explicitly testing dispersion, 25 reported using the DHARMA R
 45 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 generalized 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.

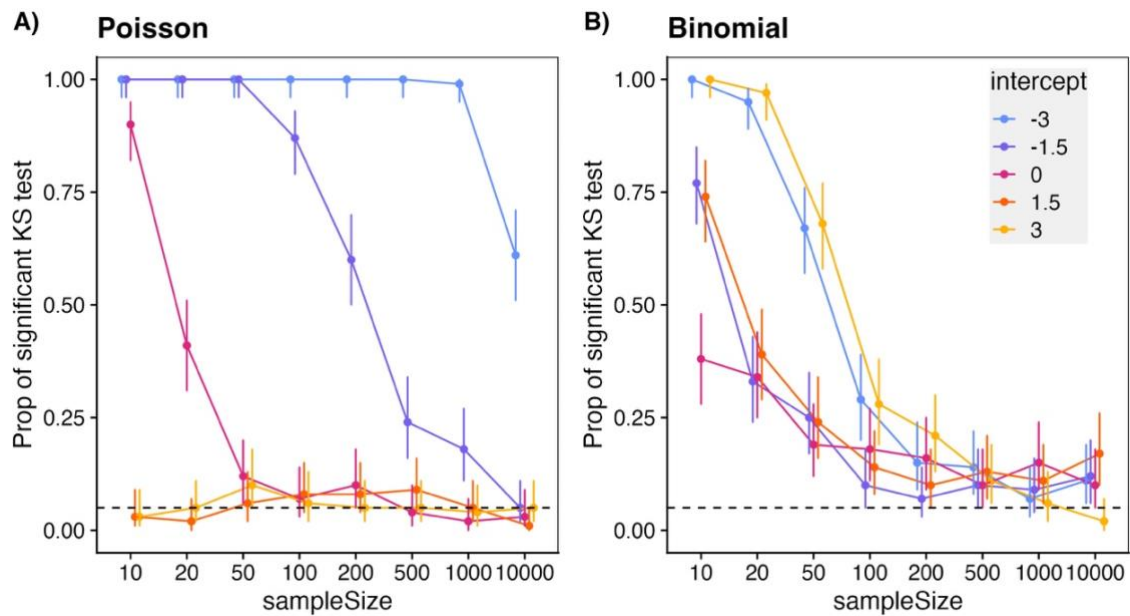
64 **S2. Pearson statistics and Chi-squared distribution**

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

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

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

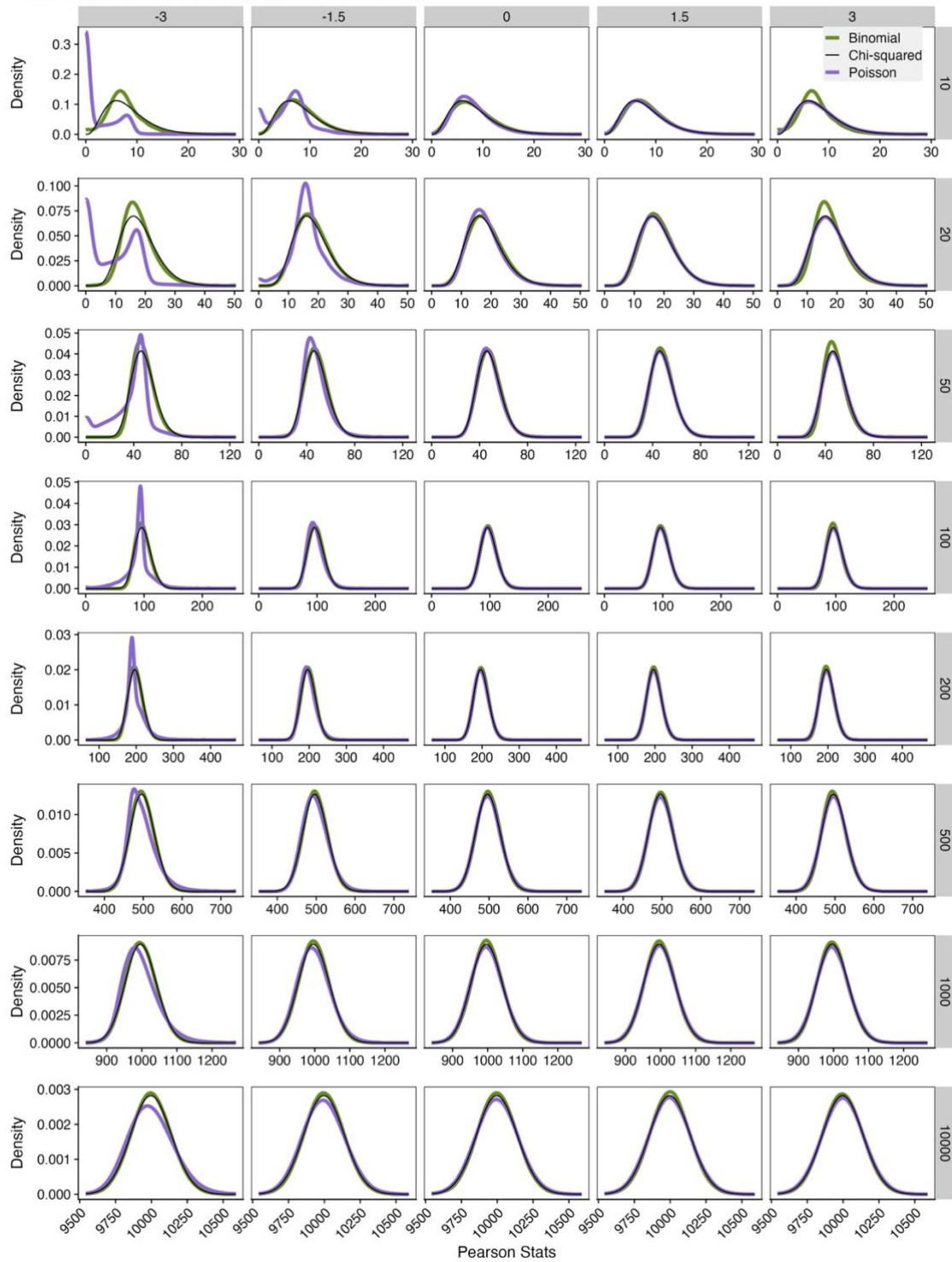
88 sample sizes (10, 20, 50) (Figure S2.1B). The proportion of significant KS tests
 89 decreased with sample size, but did not reach the nominal value of 0.05, even for very
 90 large sample sizes and intermediate intercept values (-1.5, 0, 1.5).



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



99

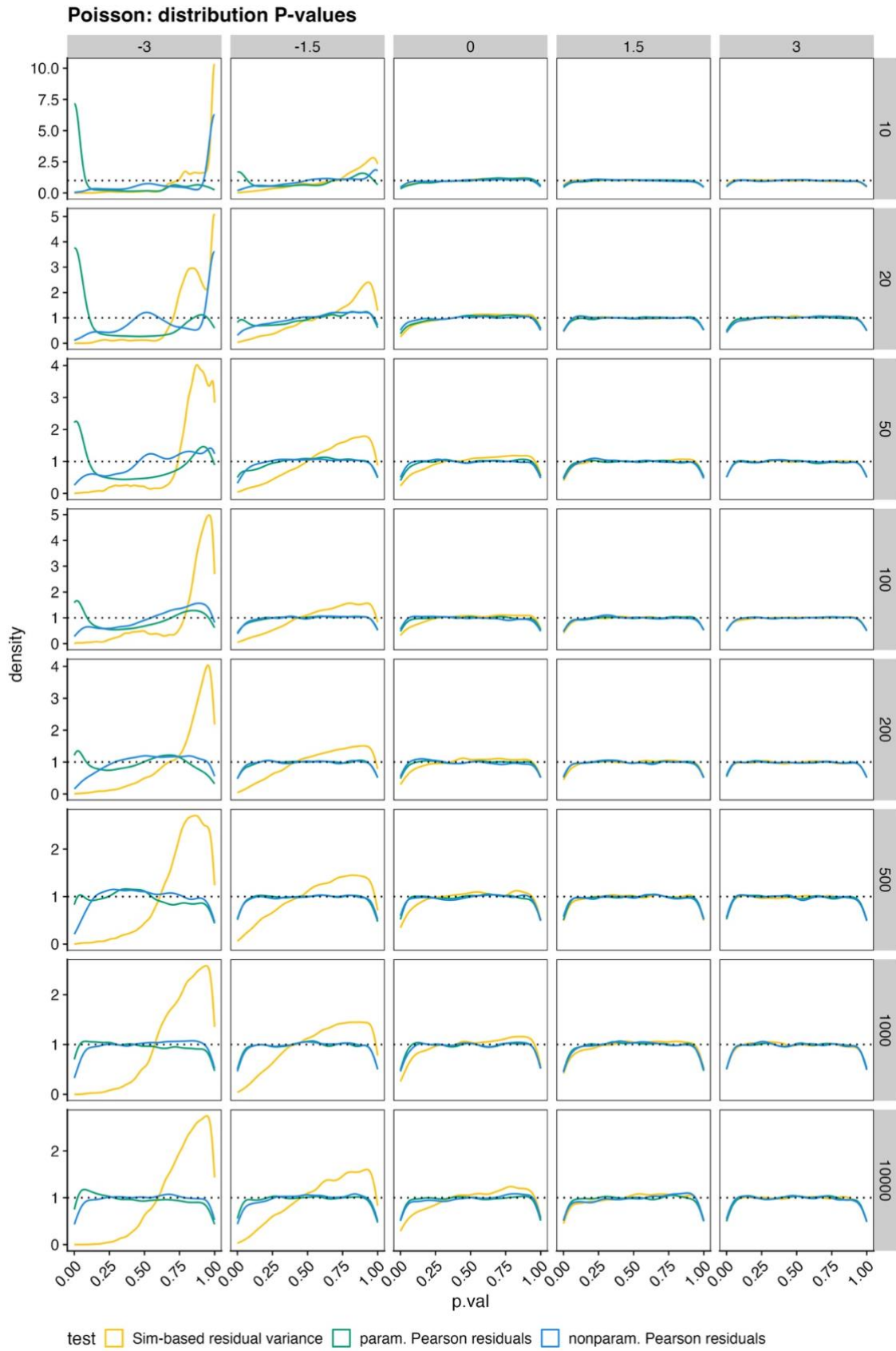
100 **Figure S2.2.** Mean Pearson statistics distribution (from 100 simulated curves) for the
 101 binomial (green) and Poisson (purple), and the Chi-square distribution in black.

102 **S3. Type I error rates for the GLMs**

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

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

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

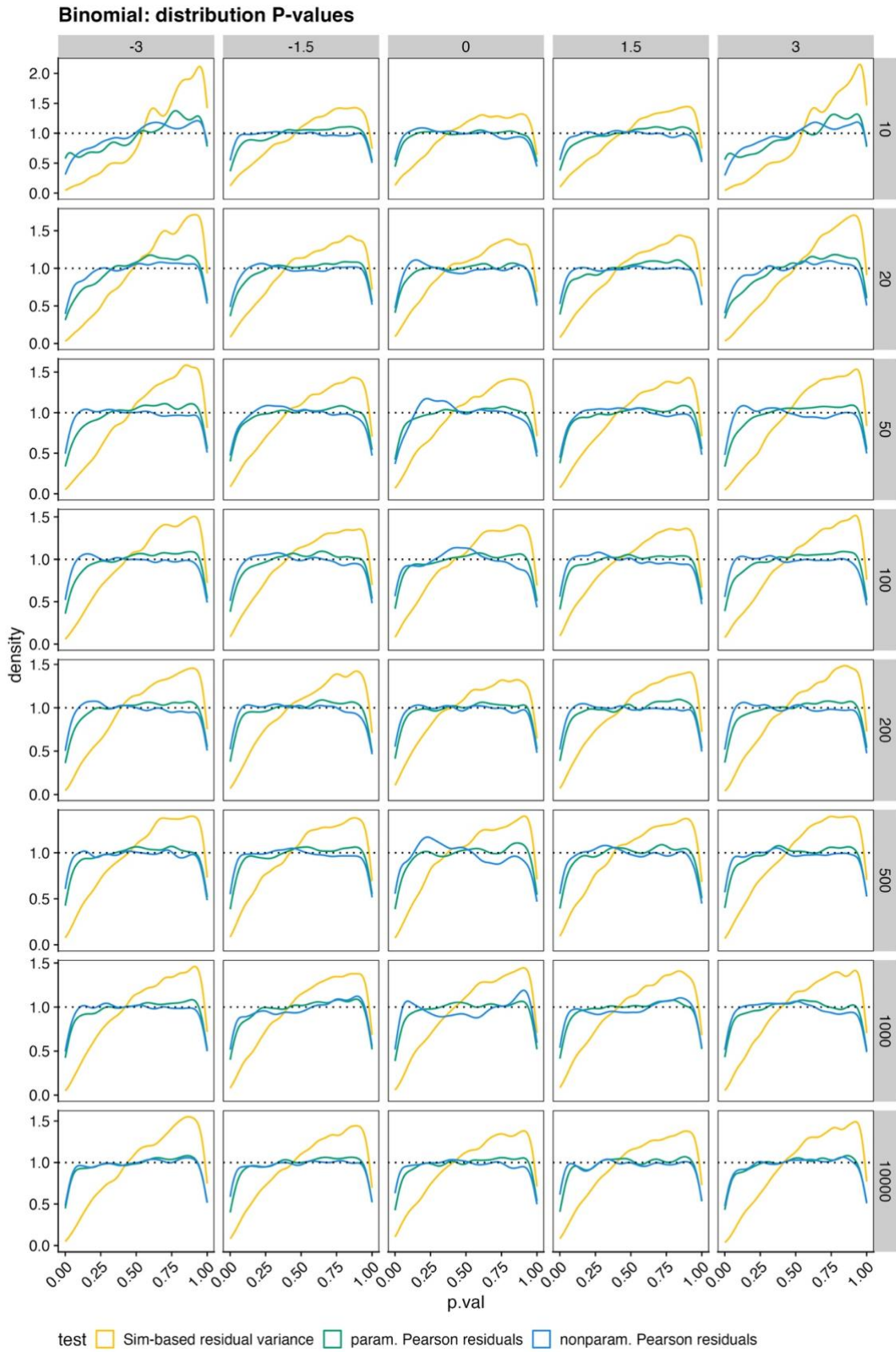


123

124

125

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).

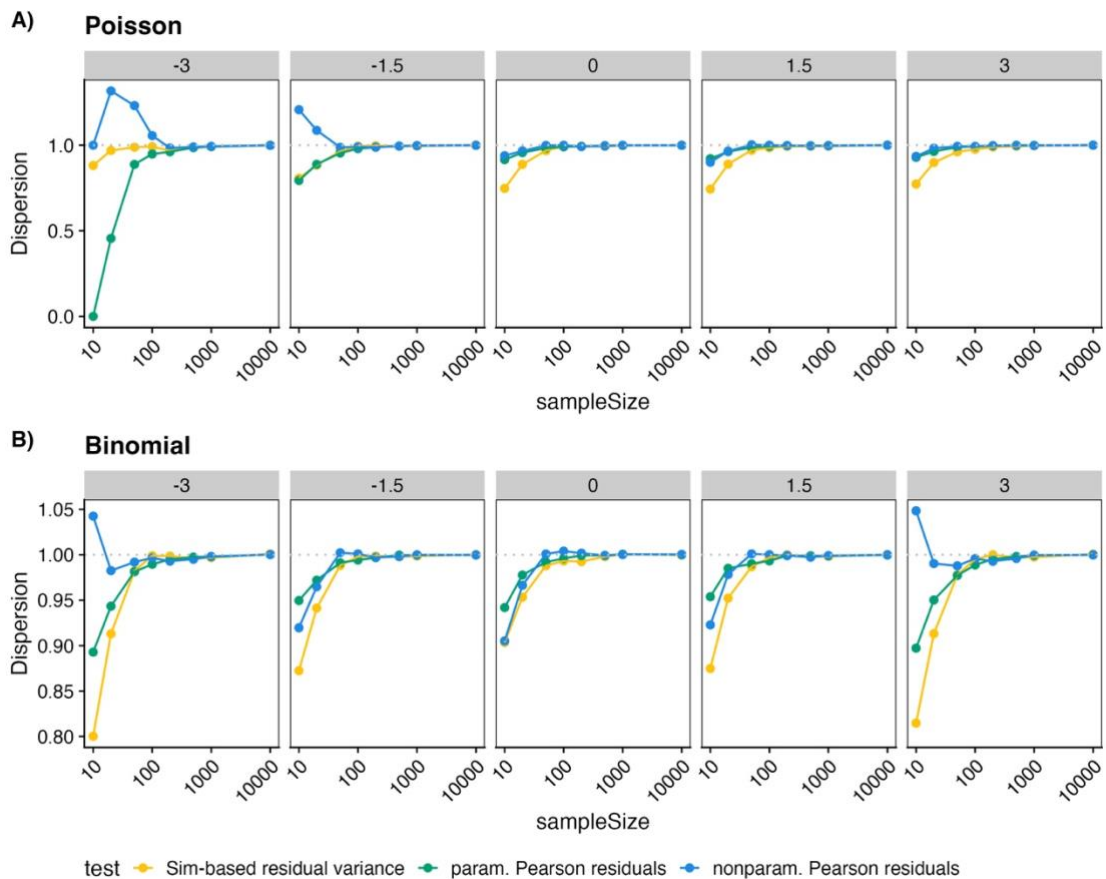


126

127 **Figure S3.2.** Distribution of p-values for the binomial GLMs for each dispersion test.
 128 10,000 simulations per simulation set (intercept x sample size).

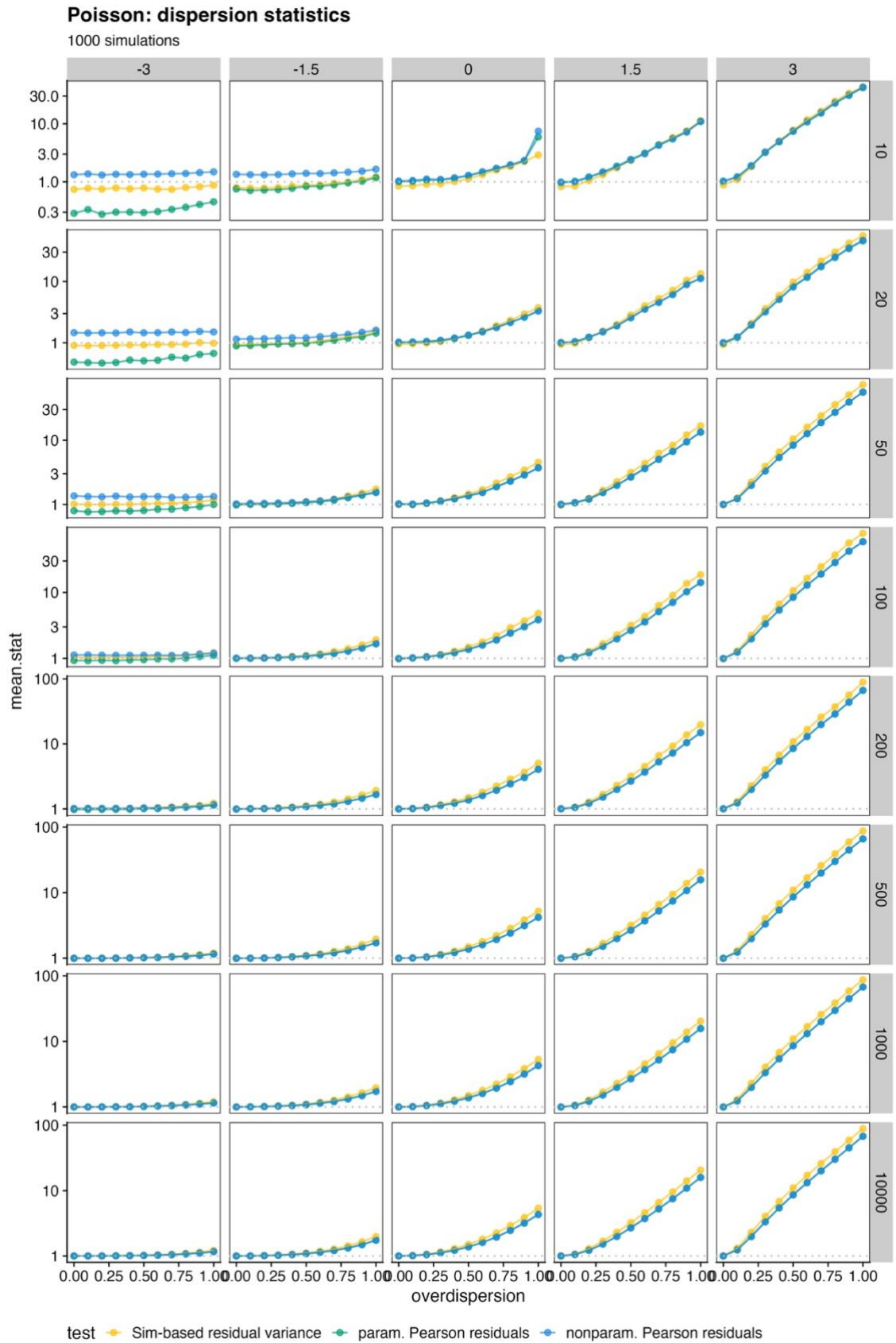
129 **S4. Dispersion statistics for GLMs**

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

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

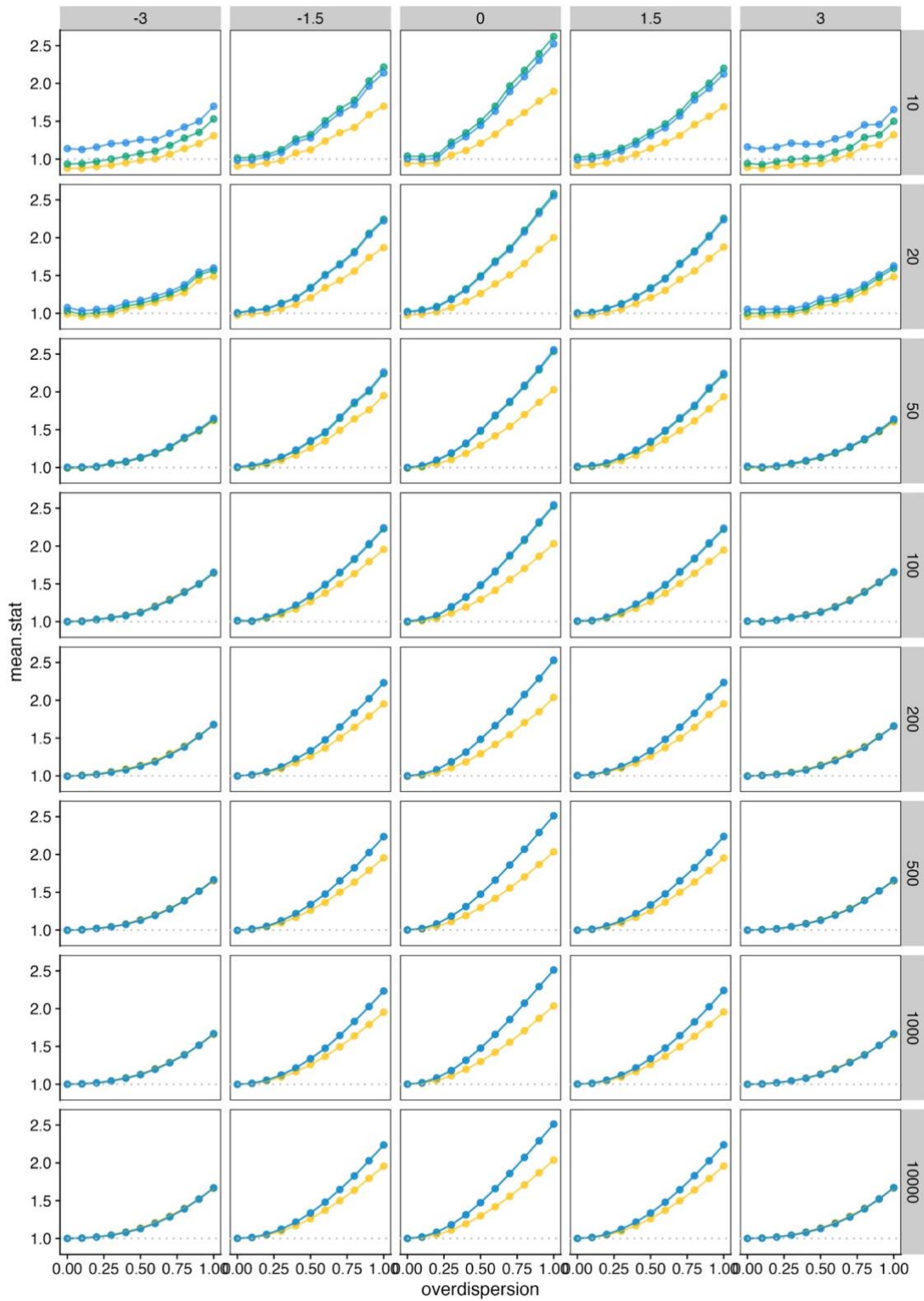


151

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

Binomial: dispersion statistics

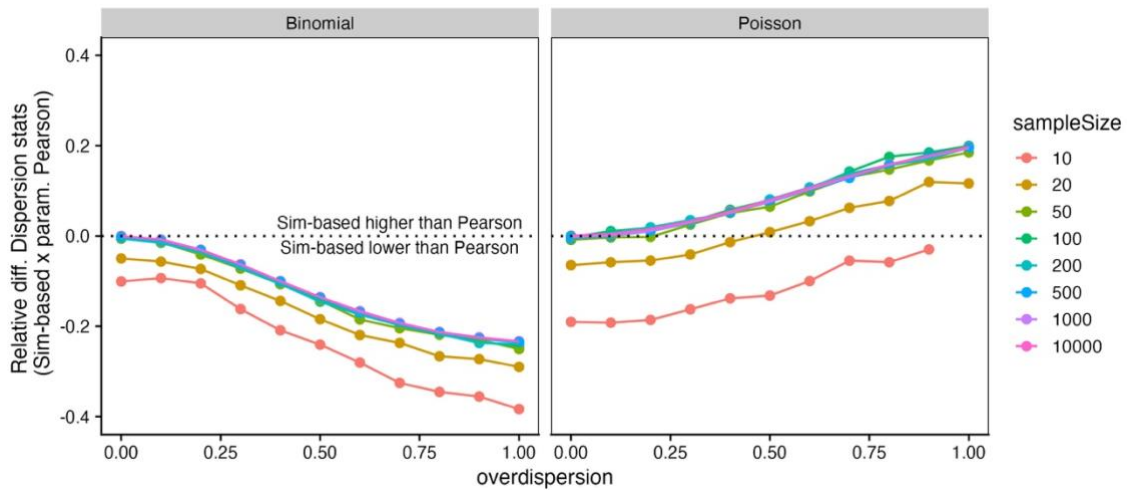
1000 sim; Ntrials=10



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

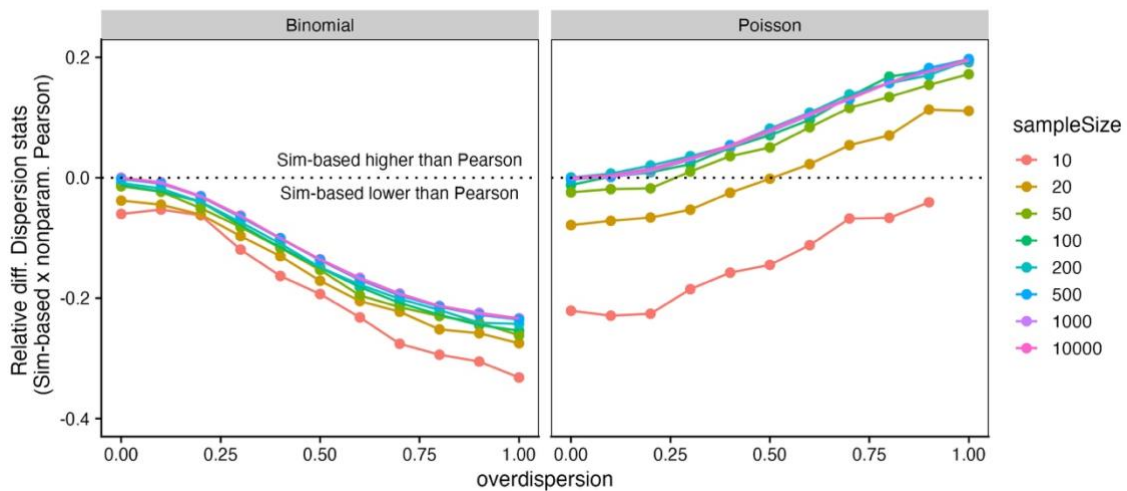
154

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



156

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

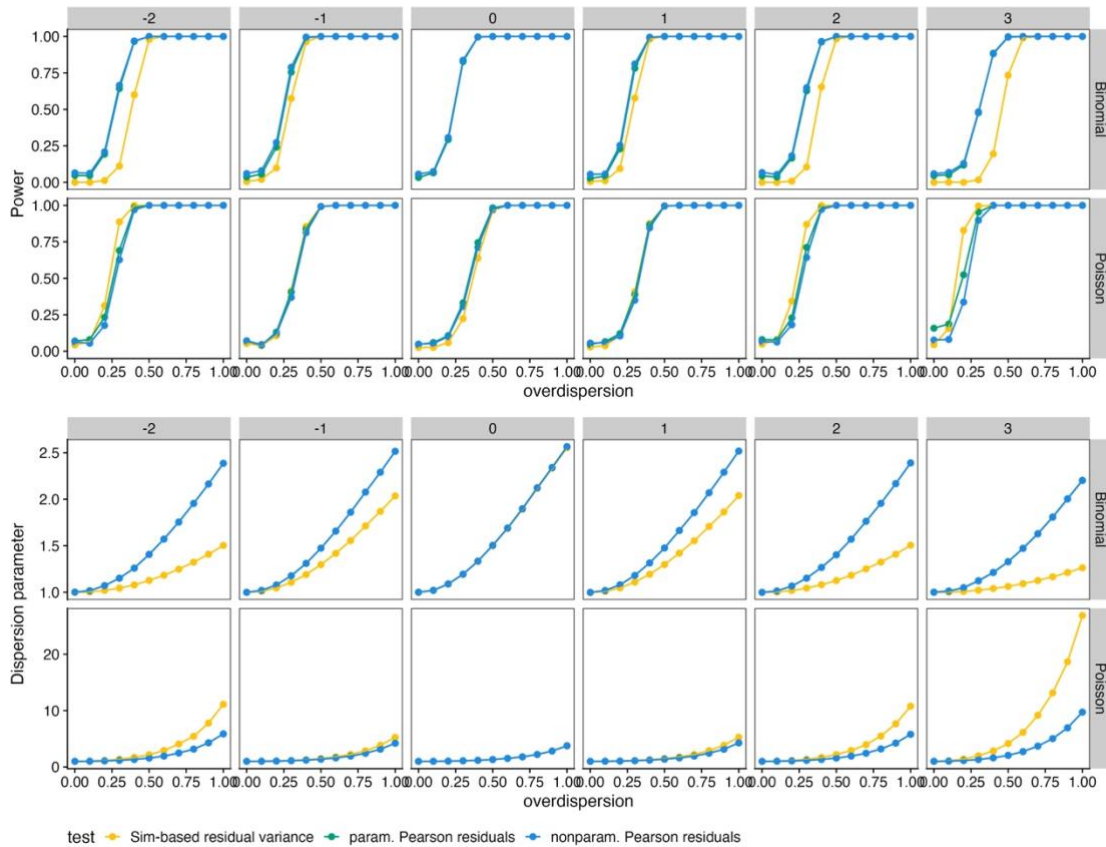


166

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

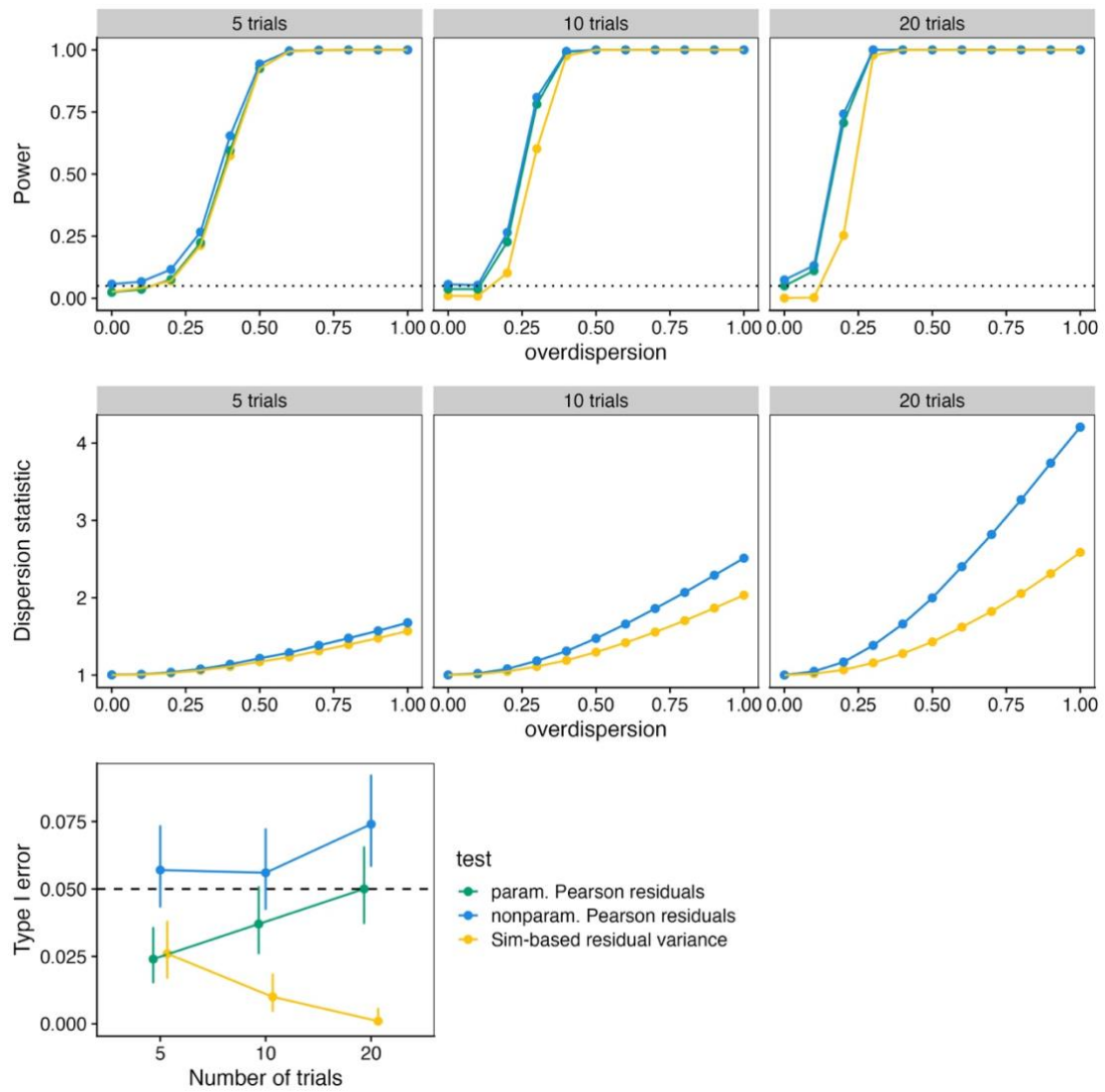
176 **S5: Expanding simulation parameters for GLMs**

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

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



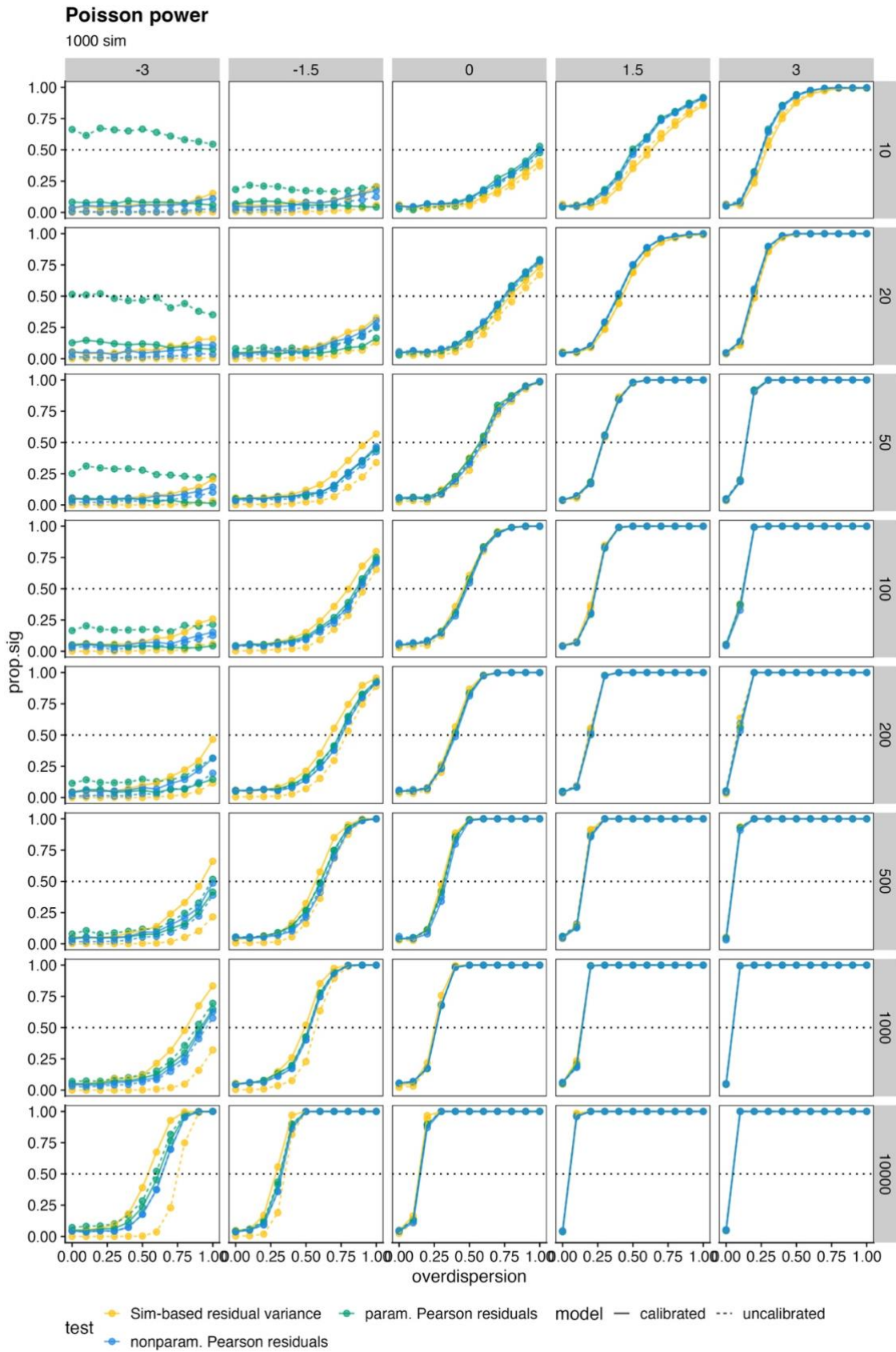
187

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

192 **S6. Power for the GLMs**

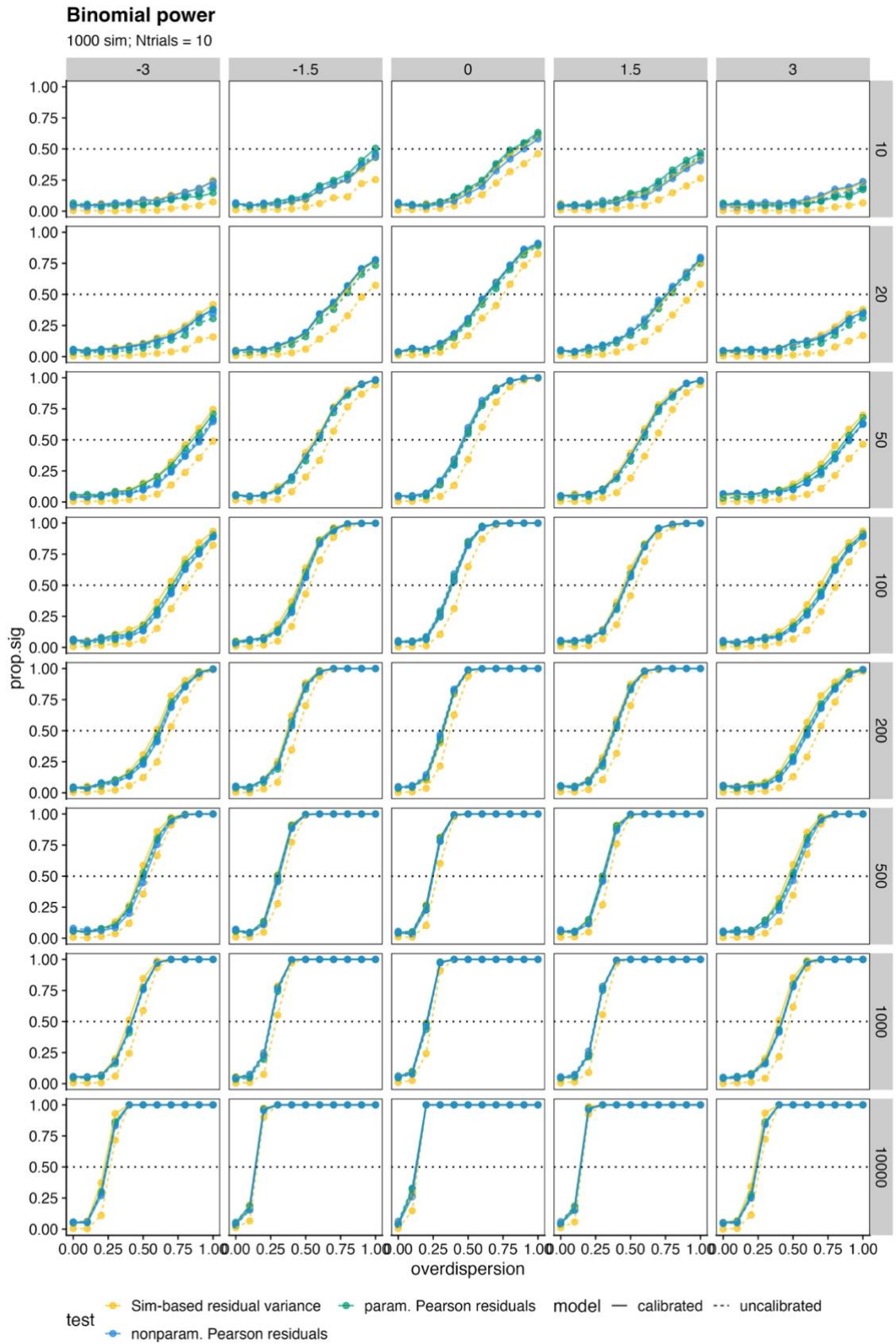
193 *Power calibration*

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



202

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



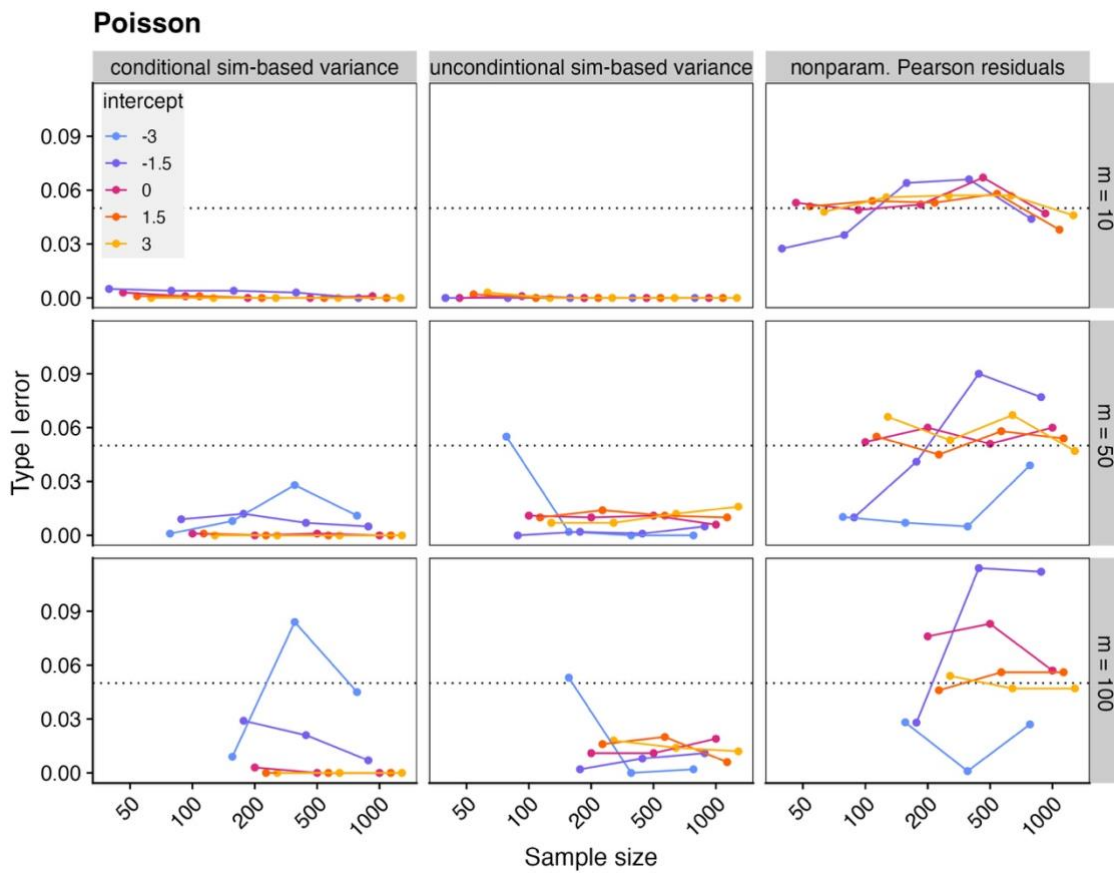
204

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

206 **S7. Additional GLMM results**

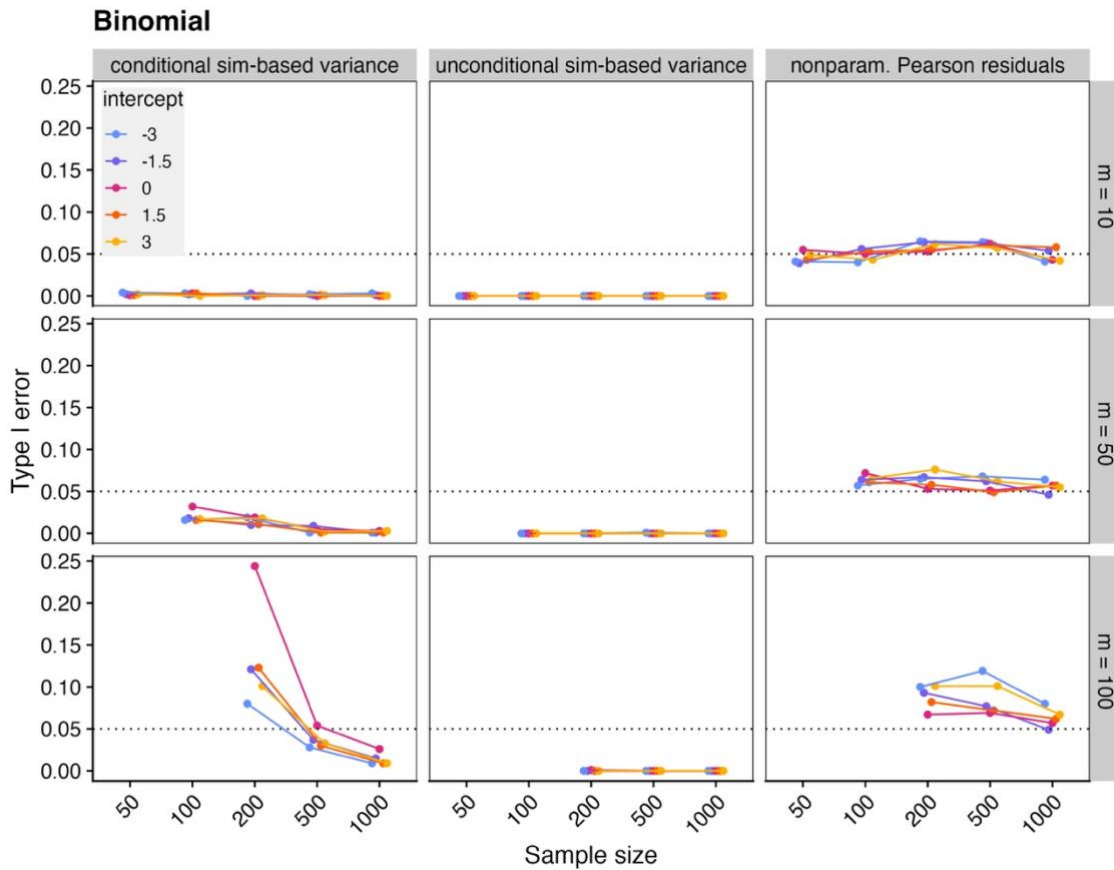
207 *Type I error rate of the alternative dispersion tests*

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

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

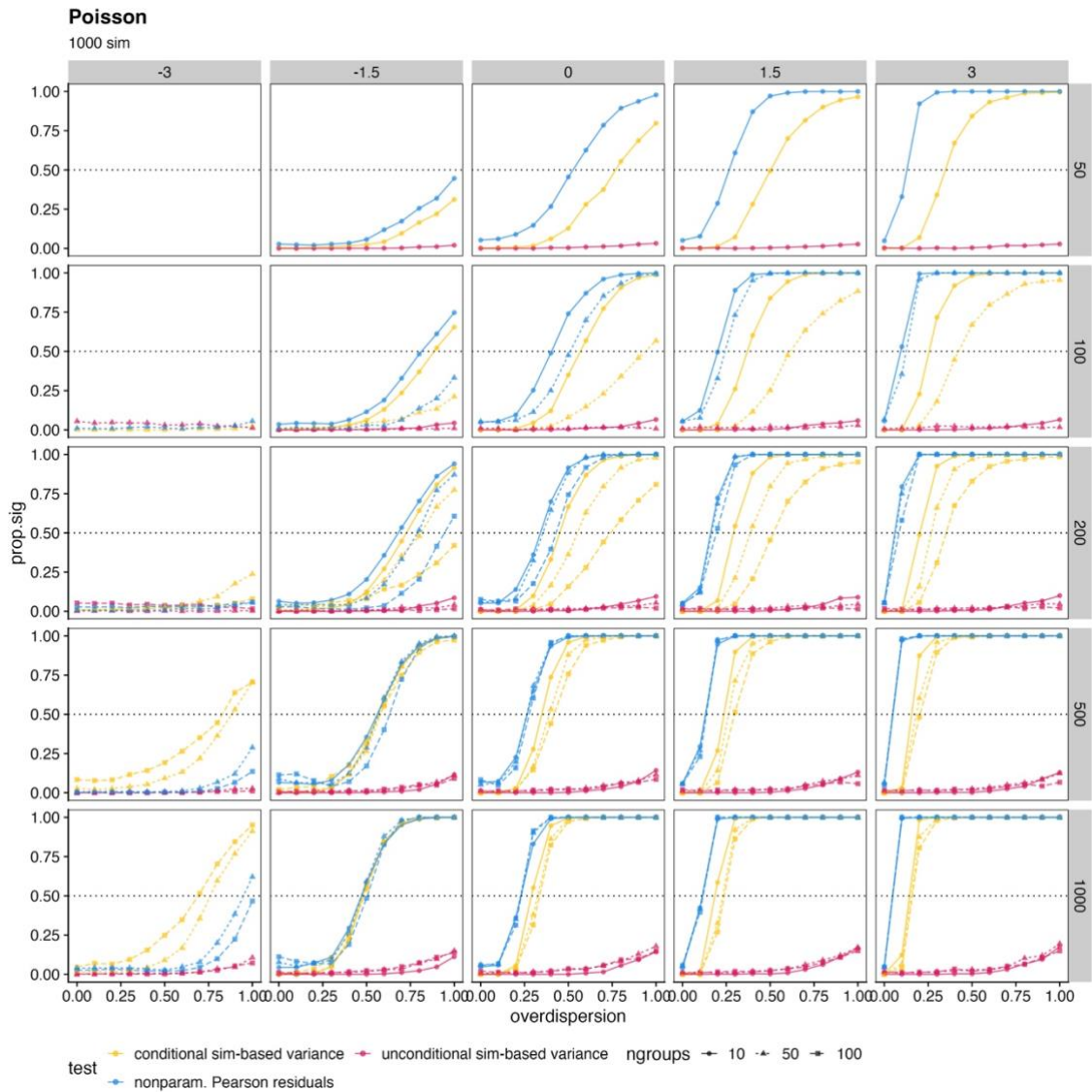


216

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

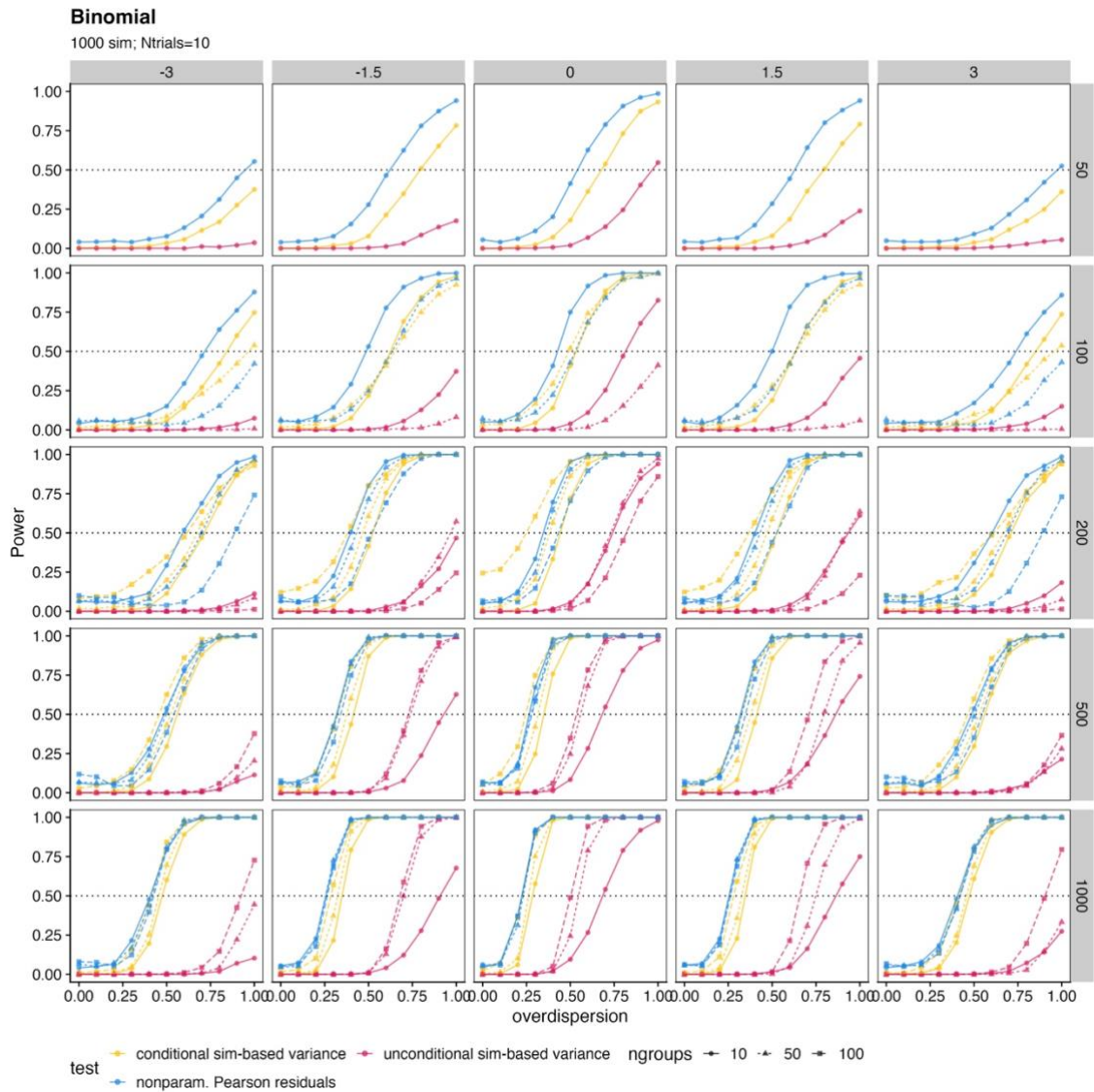
221 *Power of the alternative dispersion tests*

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



225

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

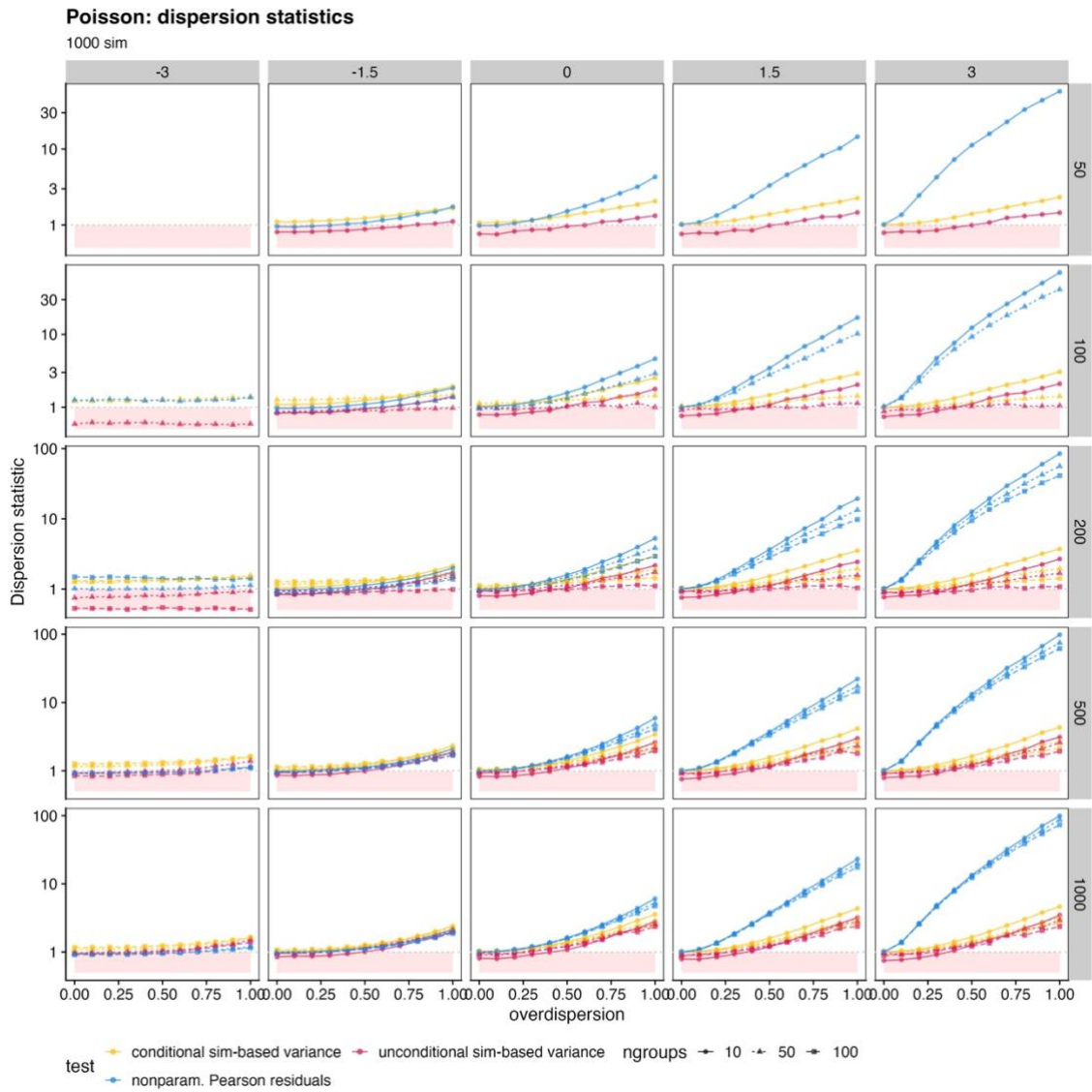


231

232 **Figure S7.4.** Power of the three alternative dispersion tests for binomial GLMMs, with
 233 different numbers of observations (rows), intercepts (columns), and number of groups
 234 for the random intercept (line types). 1000 simulations for each parameter set.

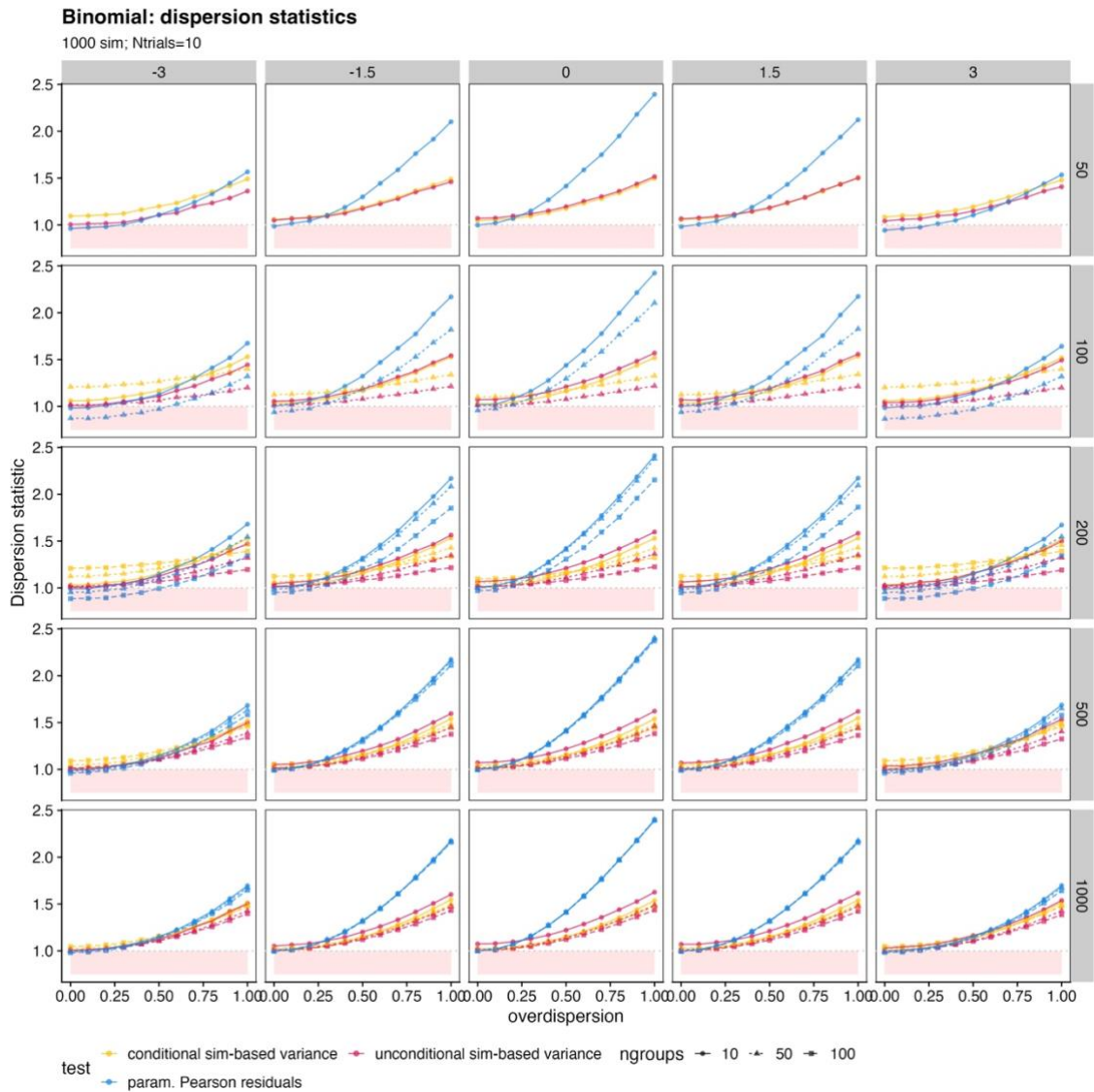
235 *Dispersion statistics of the alternative dispersion tests*

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



239

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



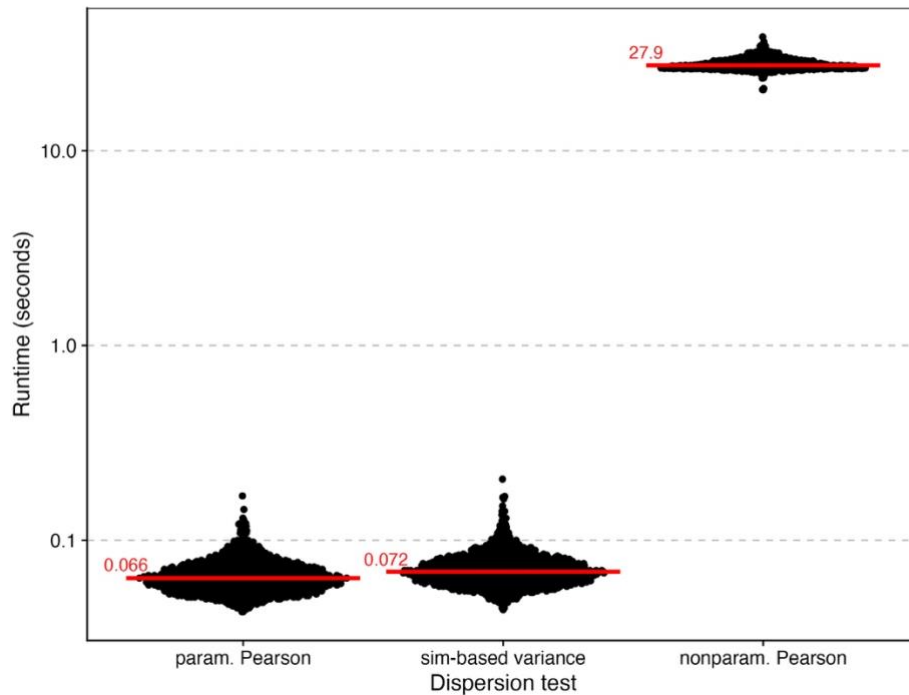
245

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

250 *Computational runtime for tests with GLMMs*

251 We computed the run time for the three tests used for GLMMs: the parametric
 252 Pearson test, the nonparametric Pearson test, and the simulation-based residual variance
 253 test with conditional simulations (Figure S7.7). We used 1,000 simulations of the
 254 Poisson GLMM as an example, with an overdispersion parameter of 0.4, an intercept of
 255 0, a sample size of 1,000, and 100 groups. There was almost no difference in

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

266 **S8: Alternative simulation-based residual variance test**

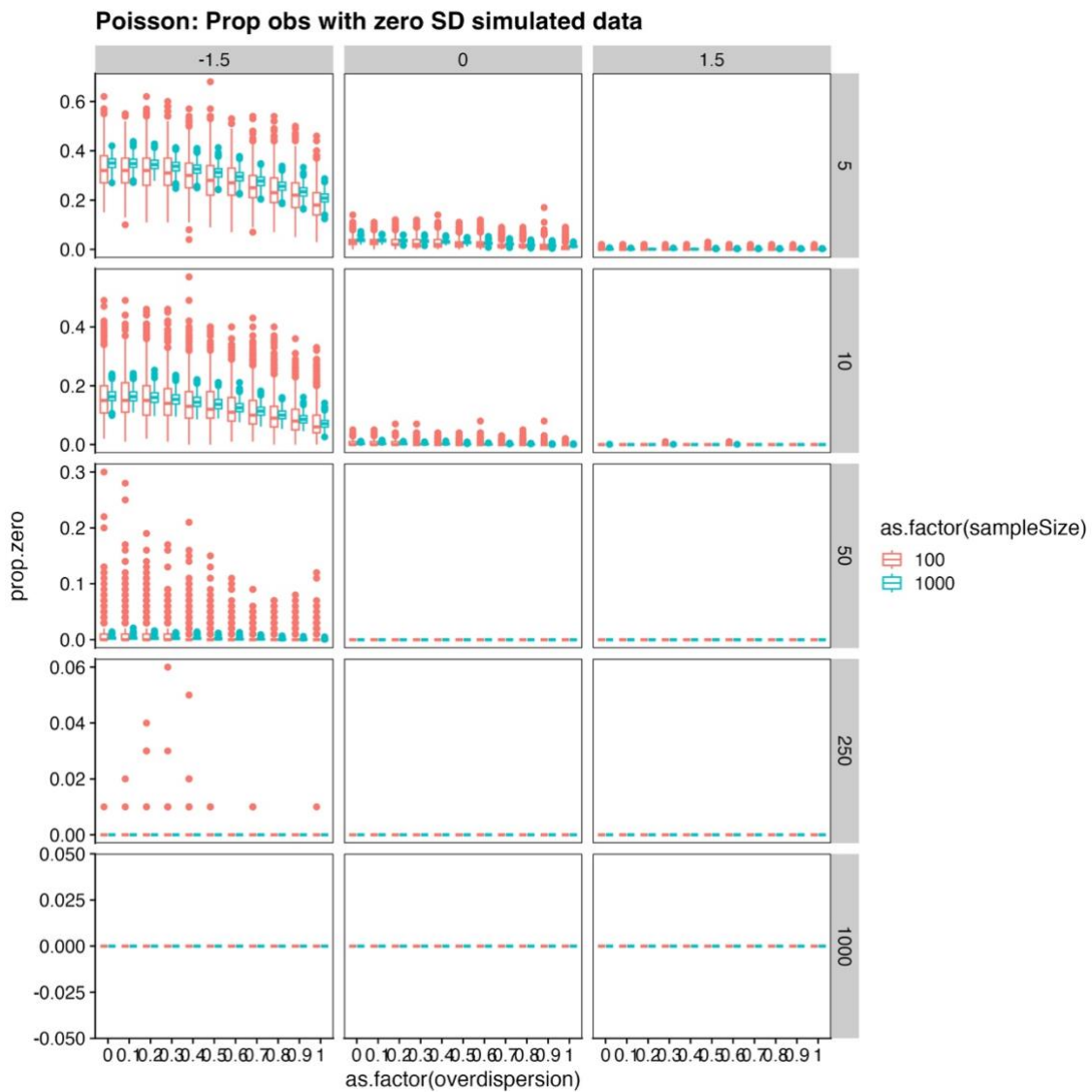
267 Another possibility for improving dispersion tests for GLMMs is to develop a
268 simulation-based approach that shows better type I, power, and a dispersion statistic that
269 could be interpreted similarly to the Pearson dispersion. To explore future possibilities,
270 we briefly considered an alternative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

$$275 \quad \textit{Approx. Pearson observed residuals: } r_i = \frac{(y_i - \hat{\mu})}{\textit{var}(y_{is})} \quad (\textit{Equation S8.1})$$

$$276 \quad \textit{Approx. Pearson simulated residuals: } r_{is} = \frac{(y_{is} - \hat{\mu})}{\textit{var}(y_{is})} \quad (\textit{Equation S8.2})$$

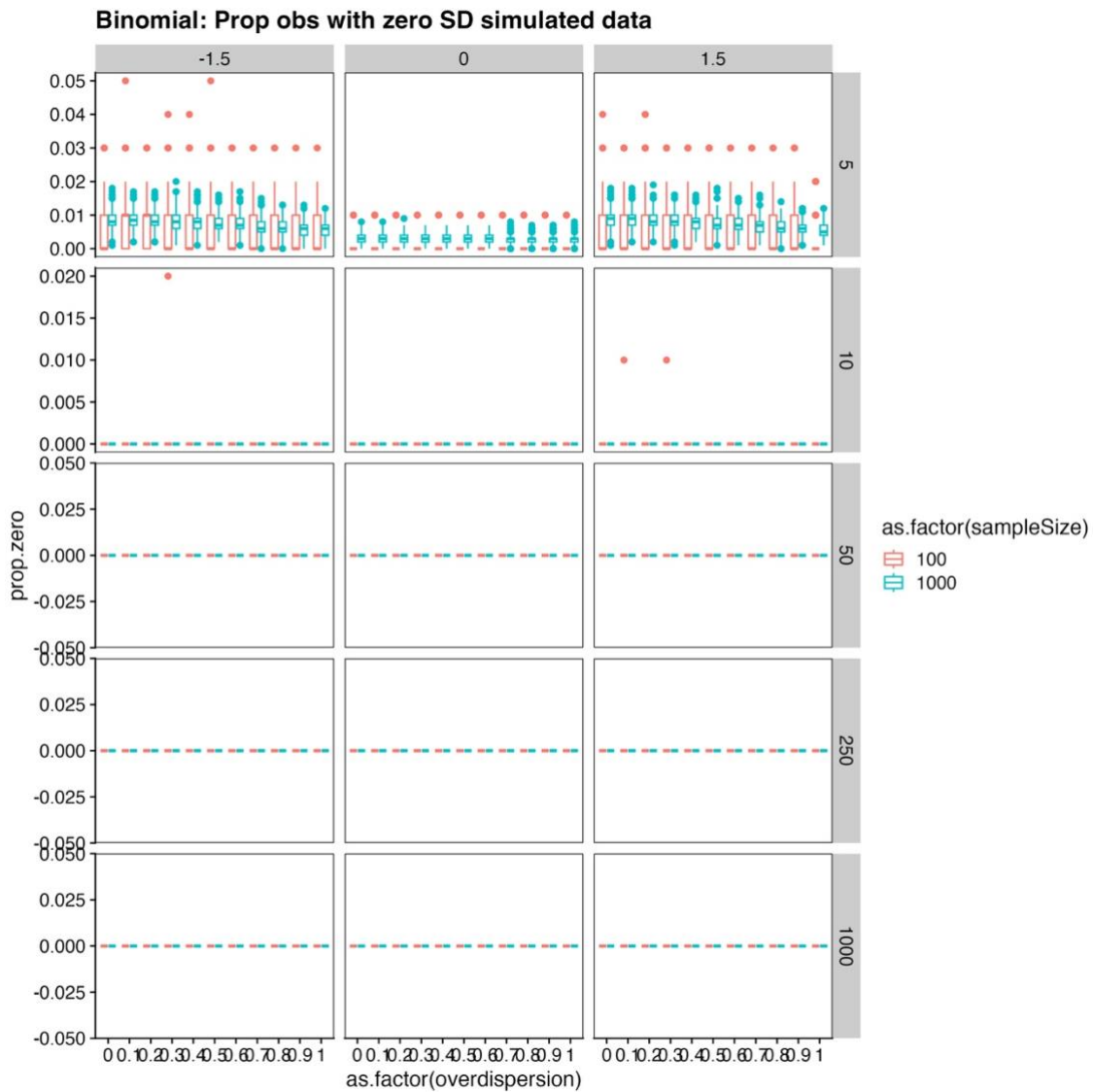
277 One obstacle with calculating the denominator of the approximate Pearson
278 residuals for each observation is that the variance depends on the number of simulations
279 and the model parameters, such as the intercept or the number of trials in the binomial
280 GLM/GLMMs. If there are too few simulations or the intercept is very small, the
281 chance of resulting in zero variance (all simulated values are the same) is higher for data
282 points with small variance. To overcome this, we first evaluated the minimum number
283 of simulations for different intercepts and sample sizes, in which all observations have
284 estimated variances that are different from zero. For all combinations of parameters, we
285 found that 1,000 simulations were sufficient to ensure that all variances in the simulated
286 observations were positive (Figures S8.1 and S8.2). However, 250 simulations (the
287 default parameter of the DHARMA package) also presented reasonable results, with the
288 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

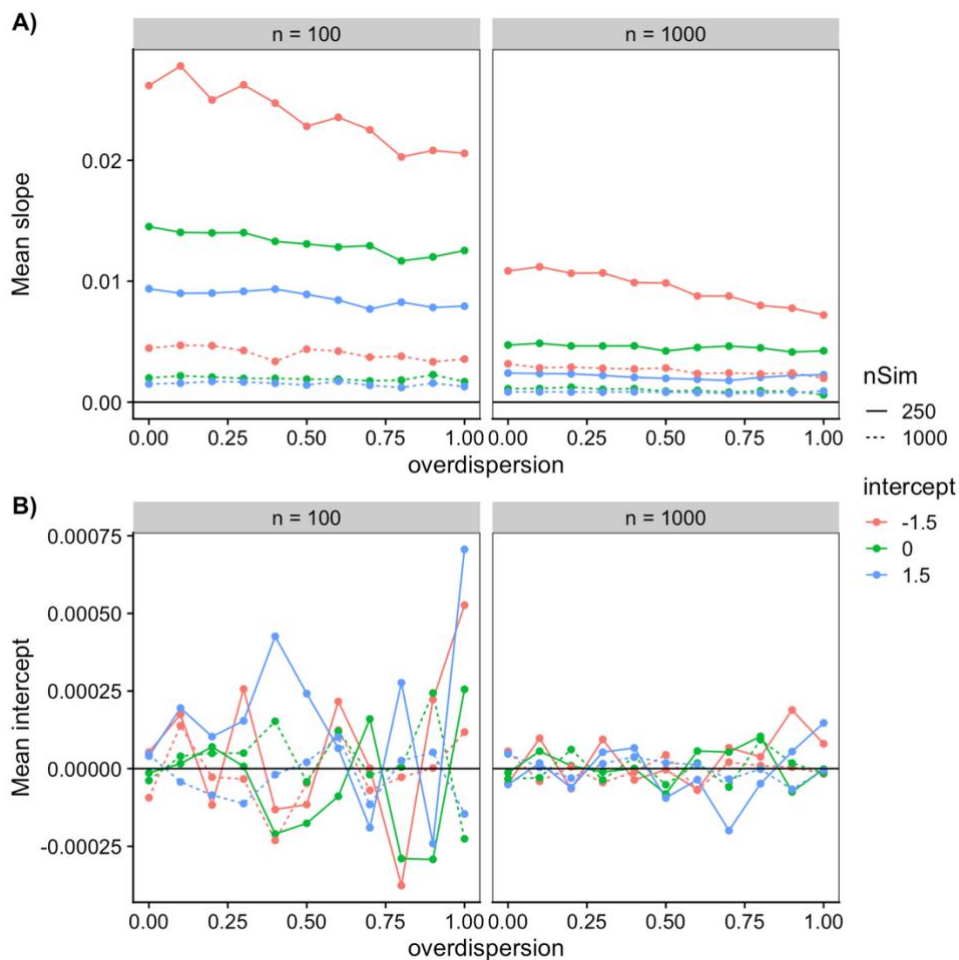
298 **Figure S8.1.** Poisson GLM: Proportion of observations with simulated zero variance in
 299 the dataset for different combinations of intercept (columns), number of simulations
 300 (rows) and sample sizes (colours).



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

306 First, we compared the approximate Pearson residuals for GLMs with the
 307 Pearson residuals by regressing the difference between them as the response variable
 308 and the Pearson residuals as the predictor for the Poisson GLMs (Figure S8.3). The
 309 intercepts for all simulation sets were nearly zero. The slope of the regression was
 310 positive and very small for the larger number of simulations and intercepts. It means

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

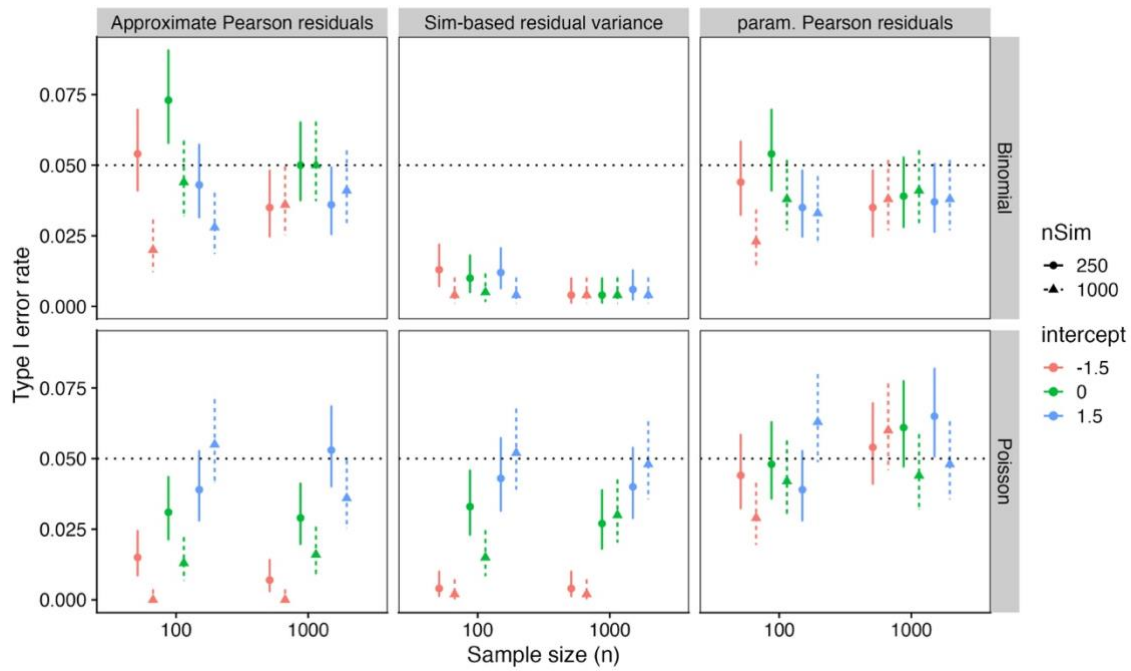


313

314 **Figure S8.3.** Mean slope (A) and intercept (B) of the regression of the difference
 315 between the Approximate Pearson residuals and Pearson residuals as response variable
 316 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).

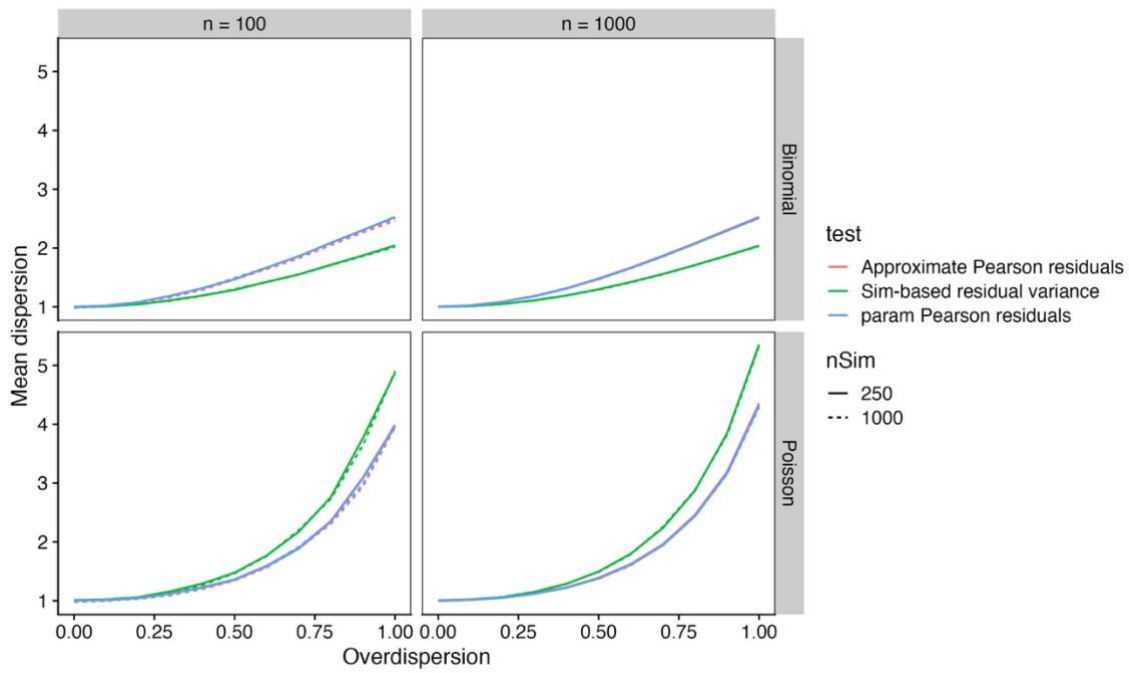
322



323

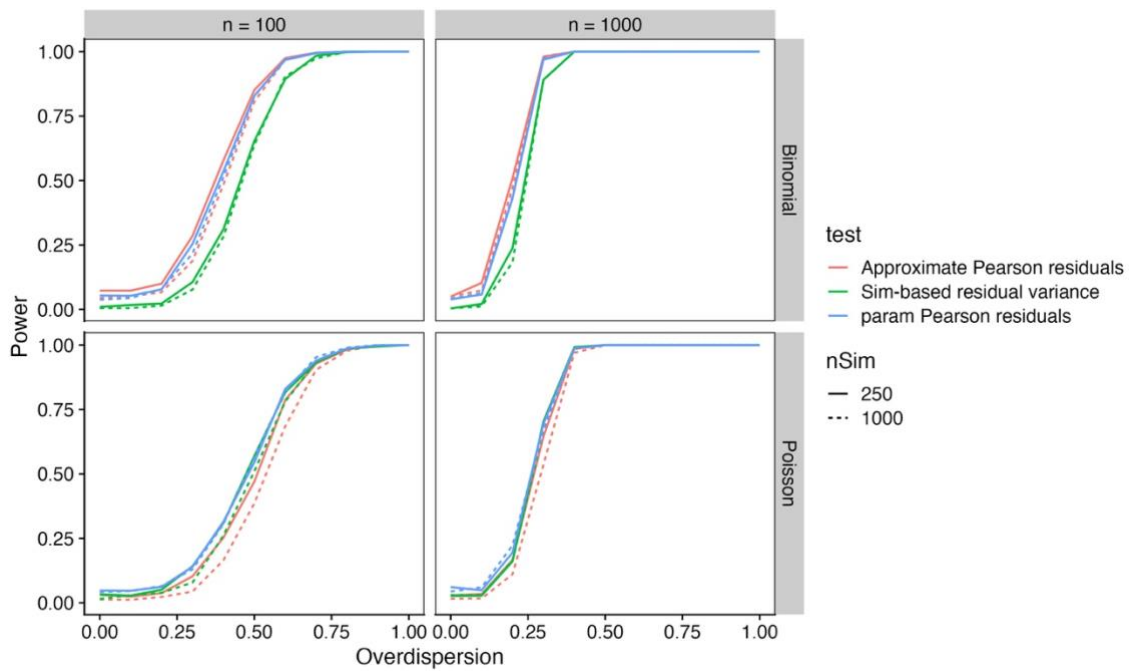
324 **Figure S8.4.** Type I error rates for GLMs comparing the parametric Pearson residuals
325 tests, the simulation-based residual variance test and the simulation-based approximate
326 Pearson test.

327 The dispersion statistics for the alternative simulation-based residual variance
328 test didn't change depending on the number of simulations and were very similar to the
329 parametric Pearson dispersion statistics for both GLMs (Figure S8.5). Power was very
330 similar among the tests for the Poisson GLM (Figure S8.6). For binomial GLMs, the
331 power of the alternative simulation-based residual test was high and similar to the
332 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

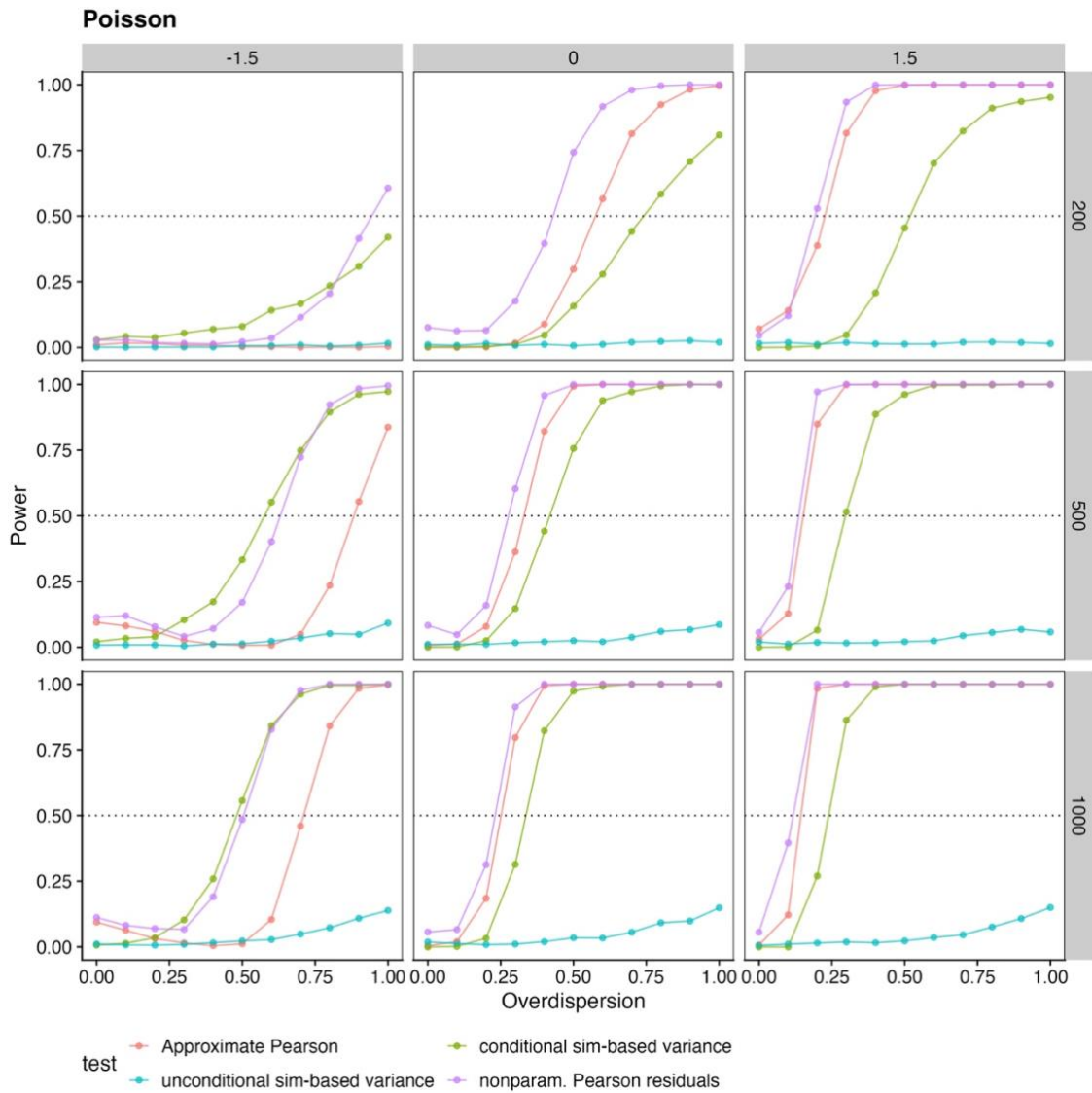
339

test fails. We compared sample sizes of 200, 500, and 1000 observations and intercepts

340

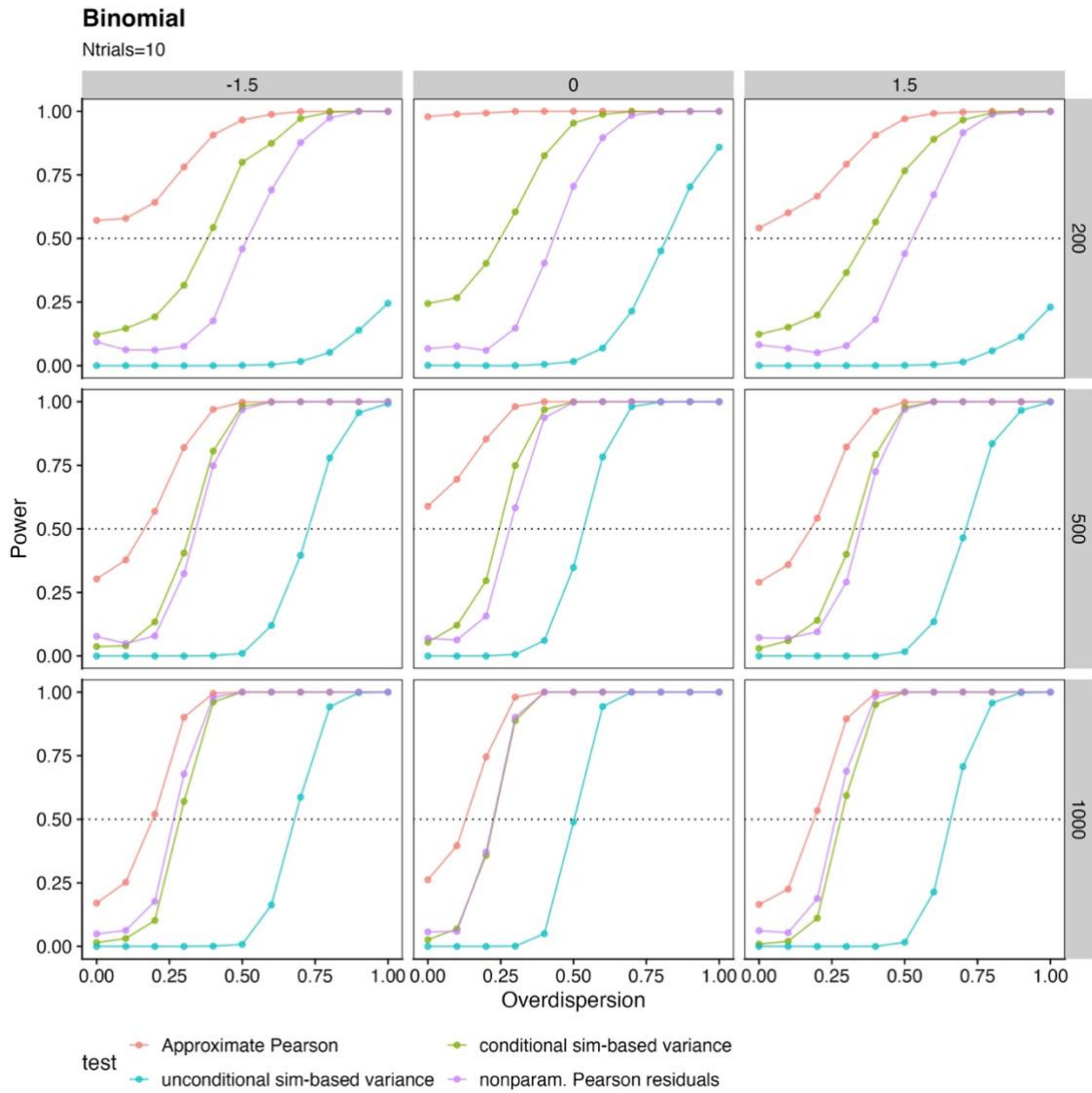
of -1.5, 0, and 1.5. We excluded simulations with zero variance in the simulated

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.



350

351 **Fig S8.8.** Power for binomial GLMMs for the alternative simulation-based test using an
 352 approximation for Pearson residuals compared with the other tests assessed in the study.
 353 1000 simulations for each parameter set: intercept (panel columns) and sample size
 354 (panel rows). The fixed parameters are slope = 1, number of groups = 100, random
 355 effects variance = 1, number of trials = 10.

356 **S9. Parametric Pearson test with approximated residual degrees of**
357 **freedom for GLMMs**

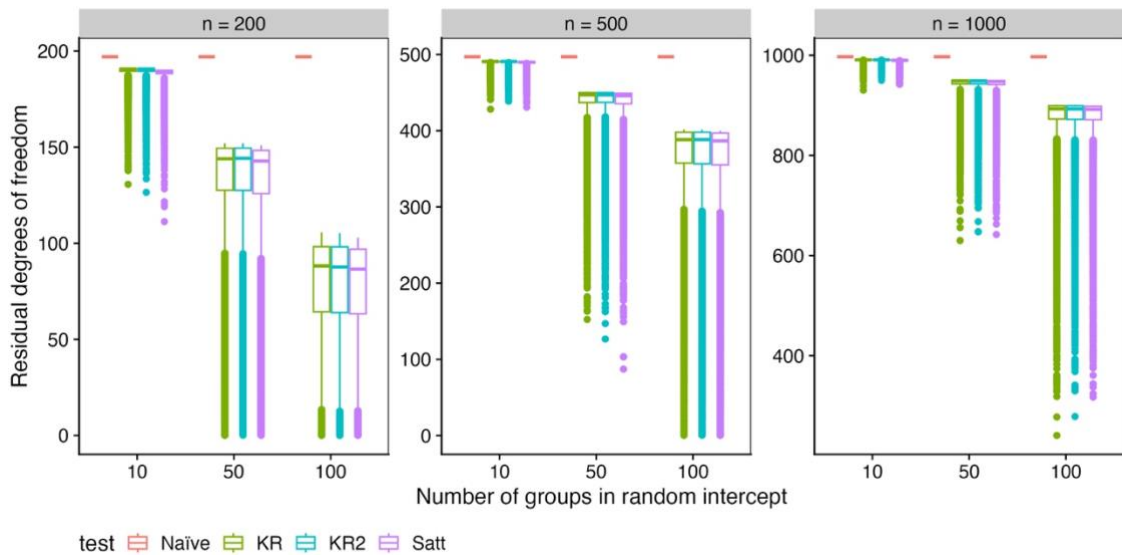
358 Degrees of freedom (*df*) are not always known for GLMMs with complex
359 hierarchical structures and limit the use of the parametric Pearson test because it
360 depends on it for evaluating overdispersion with the Chi-squared distribution.
361 Moreover, our results show that using the naïve *df* is problematic for testing dispersion
362 when you have a large number of groups in the random intercept. The two most
363 suggested methods to approximate *df* of mixed-effect models, the Satterthwaite (1946)
364 and the Kenward-Roger (Kenward & Roger 2009), were developed for LMMs to
365 account for the effect of the covariance structure on *df* and standard errors. Stroup et al.
366 (2013) suggested that the adjustment is also accurate for GLMMs. However, none of the
367 most used R packages use any correction for the degrees of freedom for GLMMs. The
368 few R packages that provide those approximations, e.g. *lmerTest* (Kuznetsova et al.,
369 2017; Kuznetsova et al., 2020) that relies on *pbkrtest* (Halekoh & Højsgaard 2014), are
370 only implemented for LMMs.

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

- 375 - The Kenward-Roger (KR) bias-corrected variance-covariance matrix for the
376 fixed effect parameters and degrees of freedom from Kenward & Roger (1997).
- 377 - The improved correction of the Kenward-Roger (KR2) returns an improved
378 correction given in Kenward & Roger (2009).
- 379 - The Satterthwaite correction (Sat) from Satterthwaite (1946).

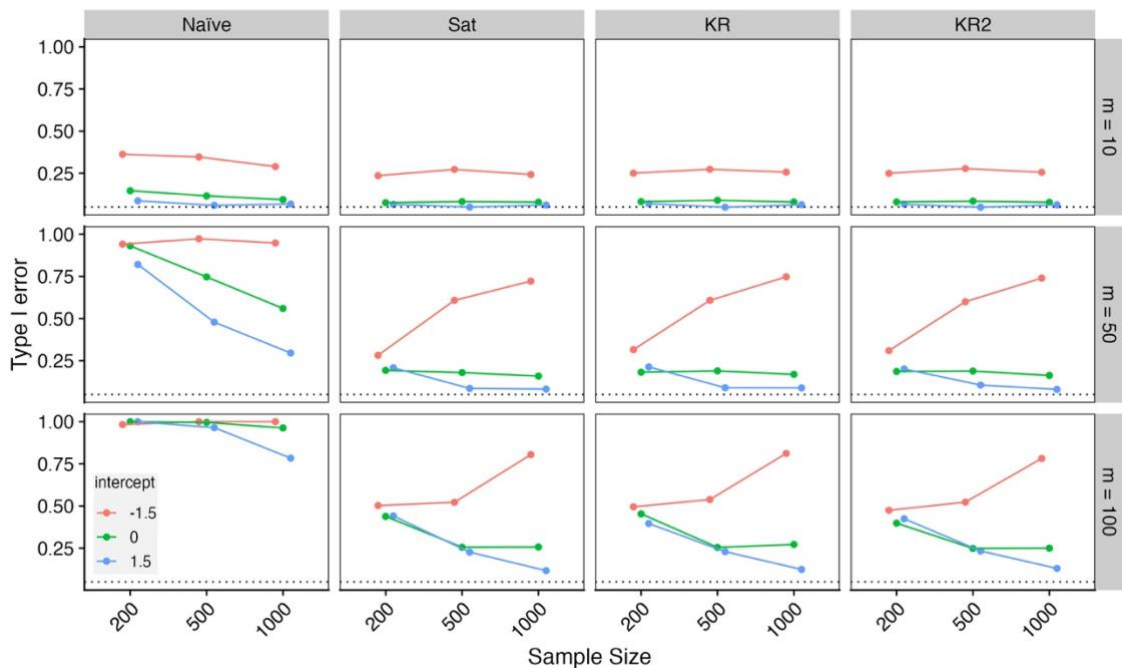
380 Our test results show that all three correction methods presented very similar
381 residual df for all simulation settings (Figure S9.1), which resulted also in very similar
382 test results (e.g., Figure S9.2 for type I error). Given the high similarity among tests for
383 the different residual 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.



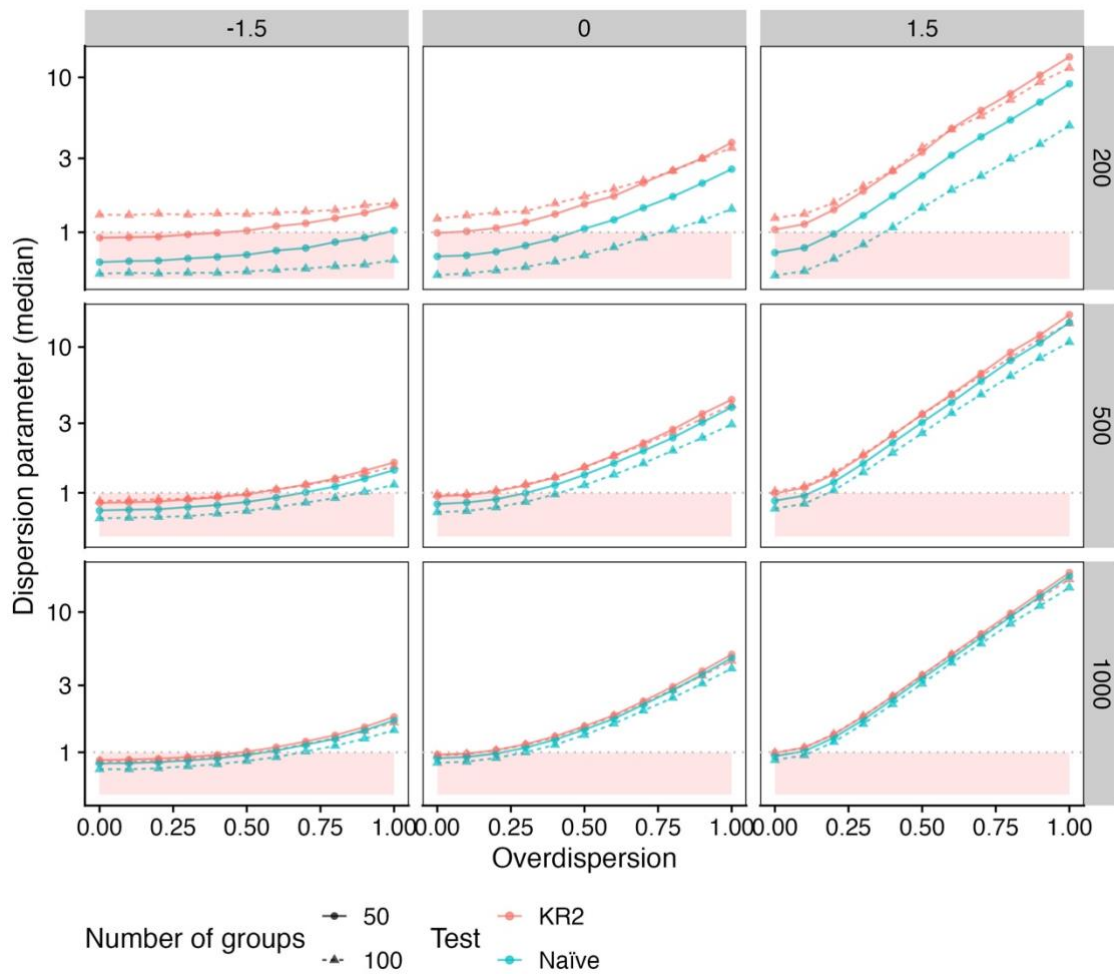
396

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



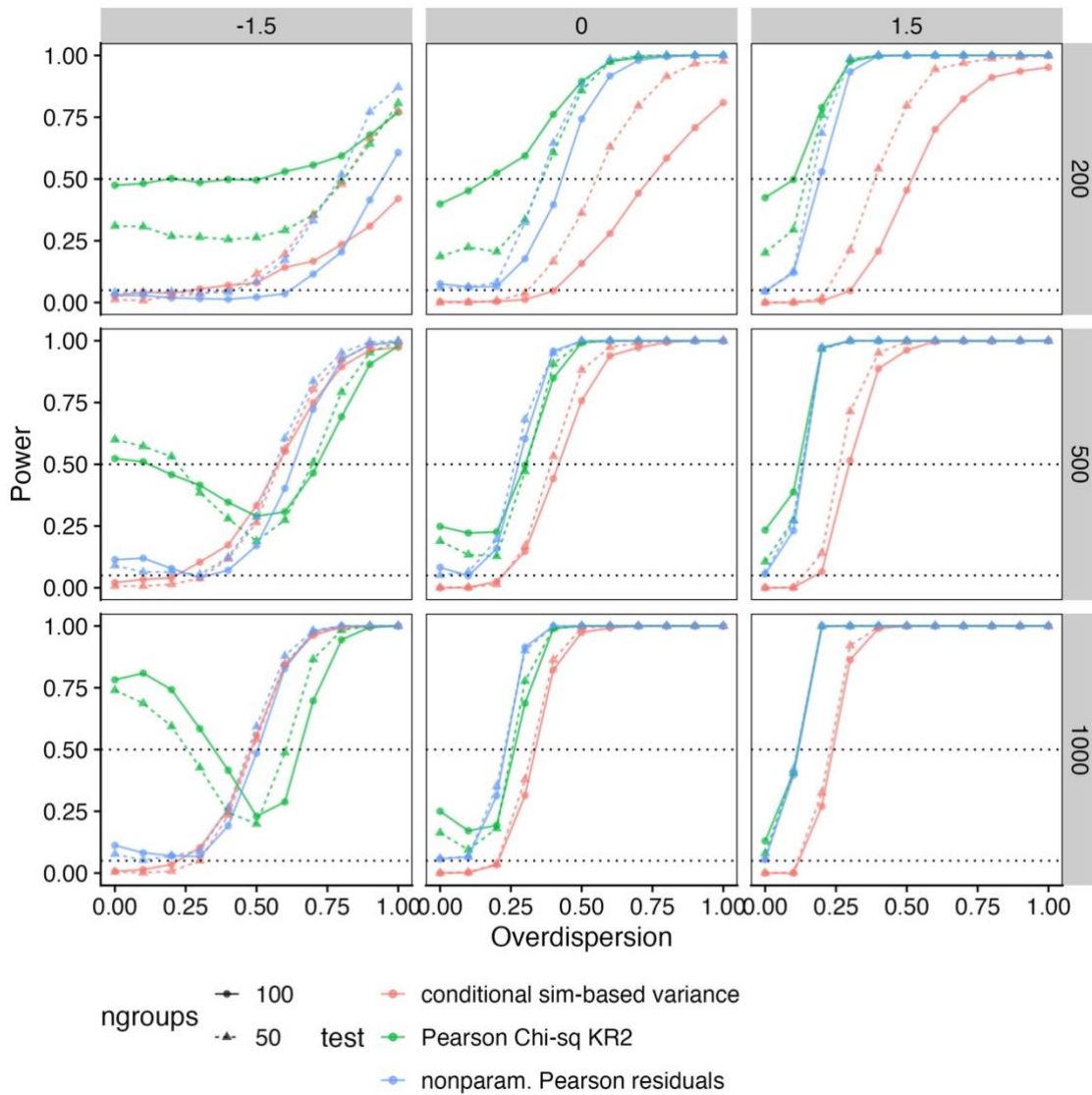
402

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



409

410 **Figure S9.3.** Dispersion parameters for the parametric Pearson test for Poisson GLMMs
 411 performed with different corrections for the residual degrees of freedom (colours),
 412 number of groups in the random intercept (linetype and shape), sample size (panel
 413 rows), and intercept (panel columns). Please refer to the main text above to relate to
 414 each applied correction. To improve clarity, we omitted the other corrections because
 415 they are too similar to each other. 1000 simulations for each parameter setting, slope =
 416 1, random intercept variance = 1.



417

418 **Figure S9.4.** Power of dispersion tests for Poisson GLMMs (colours) performed with
 419 different numbers of groups in the random intercept (linetype and shape), sample size
 420 (panel rows), and intercept (panel columns). Please refer to the main text above to relate
 421 to the applied correction for residual degrees of freedom. To improve clarity, we omitted
 422 other corrections for residual degrees of freedom because they are too similar to each
 423 other. 1000 simulations for each parameter setting, slope = 1, random intercept variance
 424 = 1.

425 **S10: Case studies detailed information**

426 *Case study 1: Redstart breeding pairs*

427 Summary results of the models applied to the Common redstart breeding pairs in
 428 Switzerland. Table S10.1 shows coefficients of models applied using a Poisson GLM
 429 (*glm* function in base R), a spatial explicit Poisson GLM (exponential spatial
 430 autocorrelation, *glmmTMB* R package) and a negative binomial GLM (*glm.nb* in MASS
 431 Rpackage). Table S10.2 shows the three dispersion tests (*testDispersion* in DHARMA)
 432 and Table S10.3 the Moran’s I residual spatial autocorrelation tests
 433 (*testSpatialAutocorrelation* in DHARMA) applied to all models.

434 **Table S10.1.** Fixed effect coefficients at the link scale (log) of the models fitted to the
 435 redstarts dataset. The effect of interest (forest cover) is highlighted in bold.

Model	Term	Estimate	SE	Statistic	P-value
Poisson	Intercept	0.137	0.105	1.314	0.189
	elevation	0.147	0.093	1.578	0.115
	elevation^2	-0.238	0.075	-3.156	0.002
	forests	-0.007	0.003	-2.213	0.027
negative binomial	Intercept	0.145	0.164	0.882	0.378
	elevation	0.106	0.138	0.77	0.441
	elevation^2	-0.231	0.101	-2.284	0.022
	forests	-0.007	0.005	-1.61	0.107
spatial Poisson	Intercept	-0.479	0.235	-2.037	0.042
	elevation	0.196	0.174	1.121	0.262
	elevation^2	-0.239	0.109	-2.203	0.028
	forests	-0.006	0.005	-1.429	0.153

436

437

438 **Table S10.2.** Residuals dispersion tests of the models fitted to the redstarts dataset. The
 439 parametric Pearson test for the spatial Poisson GLM is just for overdispersion, and the
 440 very small dispersion coefficient happens because we used a GLMM structure to model
 441 it.

Model	Test	Dispersion	P-value
Poisson	sim-based res. variance	2.47	0
	parametric Pearson	2.38	0
	nonparametric Pearson	2.39	0
negative binomial	sim-based res. variance	1.03	0.808
	parametric Pearson	1.02	0.775
	nonparametric Pearson	1.04	0.472
spatial Poisson	sim-based res. variance	0.81	0.936
	parametric Pearson	0.39	1
	nonparametric Pearson	0.99	0.984

442

443 **Table S10.3** Moran's I residual spatial autocorrelation of the models fitted to the
 444 redstart dataset.

Model	Moran's I	Expected	SD	P-value
Poisson	0.0258	-0.0029	0.0043	0
negative binomial	0.0078	-0.0029	0.0043	0.0128
spatial Poisson	-0.0003	-0.0029	0.0043	0.5407

445

446 *Case study 2: Wild and zoo-housed orangutan behavior*

447

448 **References**

- 449 Hartig, F. (2024). *DHARMA: Residual Diagnostics for Hierarchical (Multi-Level /*
450 *Mixed) Regression Models* (Version 0.4.7) [Computer software].
451 <https://CRAN.R-project.org/package=DHARMA>
- 452 Jahn, N. (2023). *europemc: R interface to the europe PubMed central restful web*
453 *service* (Version 0.4.3) [Computer software]. [https://CRAN.R-](https://CRAN.R-project.org/package=europemc)
454 [project.org/package=europemc](https://CRAN.R-project.org/package=europemc)
- 455 Laumer, I. B., Kansal, S., Van Cauwenberghe, A., Rahmaeti, T., Setia, T. M., Mundry,
456 R., Haun, D., & Schuppli, C. (2025). Wild and zoo-housed orangutans differ in
457 how they explore objects. *Scientific Reports*, *15*(1), 14853.
458 <https://doi.org/10.1038/s41598-025-97926-z>
- 459 Lüdecke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021).
460 performance: An R package for assessment, comparison and testing of statistical
461 models. *Journal of Open Source Software*, *6*(60), 3139.
462 <https://doi.org/10.21105/joss.03139>
- 463