



## 22 **Abstract**

- 23 1. Underdispersion and overdispersion are common issues when analysing  
24 ecological data with generalised linear (mixed) models (GLMs/GLMMs).  
25 Overdispersion, the phenomenon where observations spread wider than expected  
26 by the fitted model, leads to anti-conservative p-values and, thus, to inflated type  
27 I error. In contrast, underdispersion, a narrower spread of the data than expected,  
28 causes overly conservative p-values and, therefore, a reduction in power. A  
29 range of tests has been suggested to detect such dispersion problems, but there  
30 are few comparative studies of their performance across a range of models and  
31 analysis situations.
- 32 2. The goal of this study is to identify a general dispersion test for GLMs/GLMMs  
33 that is applicable across all standard distributions and random-effects structures.  
34 After an initial assessment of available tests, we selected two classes of  
35 dispersion tests as candidates: (1) parametric and nonparametric tests based on  
36 Pearson residuals and (2) simulation-based tests that compare the expected to the  
37 observed variance in the response.
- 38 3. Comparing their performance by type I error, power, and dispersion estimate,  
39 across a range of GLMs and GLMMs, we find that a nonparametric Pearson  
40 residuals test performed best across all metrics, especially for data with low  
41 incidence or count rates and/or sample sizes; however, at the cost of high  
42 computational expenses. The parametric Pearson residuals test, which is  
43 recommended in many books and guidelines, is faster and performs excellently  
44 for GLMs, but can be seriously biased towards underdispersion for GLMMs. We  
45 show that the reason for this bias, which increases with the number of random  
46 effect clusters/groups, lies in the naïve computations of the degrees of freedom

47 for the random effects. The simulation-based response variance test is slightly  
48 less powerful than the nonparametric Pearson test, but it showed overall good  
49 calibration and is much faster to compute. It offers a compromise between the  
50 strengths and weaknesses of the two Pearson-based tests.

51 4. We conclude that for GLMs, the parametric Pearson residuals test offers the best  
52 combination of speed and accuracy. For GLMMs, we recommend either the  
53 computationally demanding non-parametric Pearson residuals test or the faster,  
54 although somewhat less powerful, simulation-based response variance test.

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

## 57 **Introduction**

58 Generalised linear models (GLMs) and generalised linear mixed models (GLMMs) are  
59 the most commonly used tools for the statistical analysis of ecological data (Bolker et  
60 al., 2009; Lai et al., 2019; Touchon & McCoy, 2016). By incorporating mixed and  
61 random effect structures with a wide array of distributional assumptions (e.g., binomial,  
62 Poisson), GLMMs allow researchers to model nonnormal response variables (e.g.,  
63 counts, proportions, or presence-absence) while properly accounting for variation  
64 clustered in sampling units, sites, or study years (Bolker et al. 2009; McMahan & Diez  
65 2007). However, as for all parametric statistics, these models rely on the fact that  
66 residuals scatter around the regression mean with the specified distribution, and their  
67 inferential results can be seriously biased if these distributional assumptions are  
68 violated.

69         A particularly common and dreaded violation of distributional assumptions in  
70 generalised linear (mixed) models is overdispersion. Overdispersion refers to a higher  
71 variation in the observed data (and particularly the model residuals) than the fitted  
72 model assumes (Campbell, 2021; McCullagh & Nelder, 1989). Strong overdispersion  
73 usually appears in GLM distributions that assume a fixed mean-variance relationship,  
74 such as the Poisson model for count data (Harrison, 2014; J. M. Hilbe, 2014) or the  
75 binomial model for discrete proportions (Dunn & Smyth, 2018; Harrison, 2015). For  
76 example, a Poisson process assumes that we count randomly distributed points in space,  
77 but when observations are subject to spatial/temporal clustering and/or imperfect  
78 detection (Rhodes, 2015), we typically find higher dispersion than expected from a  
79 Poisson distribution. Alternatively, overdispersion may also arise from misfit, for  
80 example, by failing to include important predictors and interactions or by specifying the  
81 incorrect link function (J. M. Hilbe, 2011).

82           Overdispersion is a major concern in practical data analyses because it can have  
83   substantial anti-conservative effects on p-values, confidence intervals, and all other  
84   goodness-of-fit and precision metrics (Fig. 1, see also Rhodes, 2015). Anti-conservative  
85   means that p-values and confidence intervals are too small, leading to associated  
86   inflated false positive results (type I errors). In practice, we have encountered analyses  
87   where an overdispersed model had very small and significant p-values ( $<0.001$ ) that  
88   became nonsignificant after changing to a GLM with more appropriate dispersion (see  
89   also example in Fig. 1).

90           The counterpart to overdispersion is underdispersion, where the variation in the  
91   observed data (and, thus, model residuals) is lower than assumed by the fitted model.  
92   Reasons for underdispersion can again be that the data-generating process differs from  
93   what is assumed by the model (Lynch et al., 2014). However, in practice, it is often the  
94   result of model overfitting, i.e., having a too complex model that overfits the data.  
95   Underdispersion is somewhat less discussed in the literature, both because it is less  
96   frequent, but also because it leads to over-conservative model metrics (Fig. 1). This may  
97   seem less problematic as it does not lead to reporting “wrong” effects, but  
98   underdispersion reduce overall power and thus increase type II error. Therefore,  
99   accurate statistical inference demands that we identify and adequately deal with both  
100   underdispersion and overdispersion to minimise the risk of wrong inference.

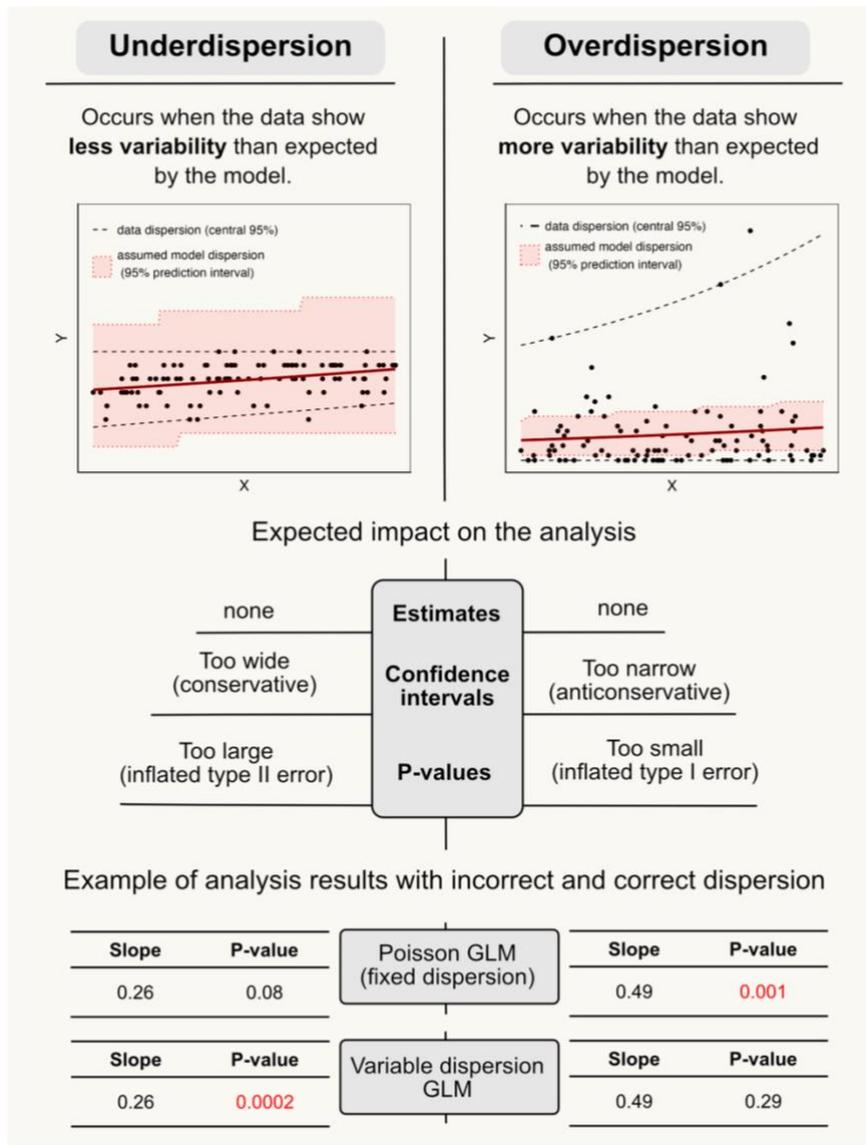
101           Due to the central importance of dispersion for all statistical indicators,  
102   statisticians have pondered how to detect and address dispersion problems since the  
103   early days of modern statistics (see Quine & Seneta (1987) and Xekalaki (2014) for a  
104   historical perspective). The first attempts to describe the phenomenon date back at least  
105   to the end of the 19<sup>th</sup> century, likely with Lexis’s ratio (Lexis 1879, apud Xekalaki,  
106   2014) for binomial clustered data, where  $Q$  is the ratio of the between-clusters variance

107 to the total variance (Xekalaki, 2014). Bortkiewicz later (1898) coined the term  
108 “divergence coefficient” ( $Q^2$ ), which is the variance divided by the mean of the sample,  
109 as a test statistic for the Poisson model (Quine & Seneta, 1987). William Gosset, the  
110 inventor of the t-test, also considered the problem of dispersion in the Poisson model  
111 (Student, 1919). Since then, a large variety of approaches have been proposed and  
112 discussed to deal with the “dispersion problem”, ranging from (1) comparing models  
113 with or without free dispersion parameters through likelihood ratio test, such as Poisson  
114 and negative binomial (e.g. Yang et al., 2007), (2) designing specific hypothesis tests for  
115 the “extra” variation (e.g. Fisher, 1950), such as score tests (Dean, 1992; Dean &  
116 Lawless, 1989; Lawless, 1987), (3) using goodness-of-fit tests, such as tests on Pearson  
117 or Deviance residuals (Dunn & Smyth, 2018; McCullagh, 1985) (although the  
118 distinction between categories (2) and (3) can be blurry, see (Collings & Margolin,  
119 1985; Dean, 1992; Dean & Lawless, 1989) or (4) using simulation-based non-  
120 parametric tests to compare observed and predicted variance of the response data  
121 (Hartig, 2024).

122         Somewhat confusing for the practical data analyst, however, many of these  
123 approaches have been designed and tested only in very specific scenarios (e.g. only for  
124 a Poisson GLM), and there is a surprising lack of systematic evaluation of these tests  
125 and strategies across a range of more complex GLMMs. Moreover, a quick review of  
126 current methods available in the R environment (R Core Team, 2024) revealed that  
127 existing dispersion tests are scattered across different packages (Table 1), and most of  
128 these only work for a restricted set of models. All this makes it challenging to decide  
129 which test should be used in an applied data analysis.

130         The goals of this study are: (1) to review and order the diversity of dispersion  
131 tests for GLMs and GLMMs, and (2) to identify tests that can reliably work across a

132 range of models with diverse distributions and complex hierarchical structures. Based  
133 on our literature review (next section), we identified two groups of tests that appeared to  
134 be generally applicable: parametric and non-parametric tests on Pearson residuals, as  
135 well as a new simulation-based non-parametric test that directly compares observed and  
136 predicted variance of the response data. We then used simulated data to compare the  
137 performance of these tests in terms of type I error, power, and the interpretability of the  
138 dispersion statistics. Based on this, we provide recommendations on the most suitable  
139 tests for detecting over- or underdispersion, depending on model complexity and  
140 software availability (i.e., currently available packages and functions in R).



141

142 **Figure 1.** Definition, statistical consequences, and a practical analysis example of  
 143 under-/overdispersion in generalised linear (mixed) models. The top row shows  
 144 examples of a data analysis using a Poisson GLM with simulated under- and  
 145 overdispersed count data. The data points in black are contrasted to the Poisson model's  
 146 95% prediction interval (in red). Black dashed lines illustrate the data dispersion  
 147 (central 95% quantiles of the data). In the example data analysis, we present slope  
 148 estimates and p-values for the GLM Poisson model fitted to the under- and  
 149 overdispersed data above, as well as the results using more appropriate models with  
 150 correct dispersion, here a Conway-Maxwell-Poisson GLM for underdispersed data and  
 151 a negative binomial GLM for overdispersed data.

**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
<b>Likelihood Ratio Test (LRT)</b>	Compare two models with and without free dispersion parameters, for example: - Poisson and negative binomial or generalized Poisson - binomial and beta-binomial	Requires fitting two models, requires defining an alternative model.  Not a dispersion test.	pscl::odTest()	GLM Poisson -> negative binomial with MASS::glm.nb()	Jackman (2024)
			DCluster::test.nb.pois()	GLM Poisson -> negative binomial with MASS::glm.nb()	Lopez-Quílez (2005)
			anova(..., test="LRT")	Many GLM/GLMMs (Different packages have the S3 method for anova functions to perform LRT)	
			lmtree::lrtest()	Any GLM	(Zeileis & Hothorn, 2002)
<b>Score test</b>	Evaluate score of restricted dispersion parameter	Requires score calculation for specific models. R functions only for Poisson GLM.	DCluster::DeanB() DCluster::DeanB2()	GLM Poisson. Score tests based on Dean (1992)	Lopez-Quílez (2005)
			Rfast2::overdispreg.test()	GLM Poisson (own model implementation)	Papadakis et al. (2025)
	Regression-based test for overdispersion from Cameron & Trivedi (1990)	Distribution specific (Poisson-based only).	overdisp::overdisp()	GLM Poisson (own model implementation)	Cameron & Trivedi (2023)
			AER::dispersiontest()	GLM Poisson from stats::glm()	Kleiber & Zeileis (2008)
<b>Standard. residuals dispersion</b>	A goodness-of-fit test to evaluate residual dispersion, e.g. via sum of Pearson residuals.	<b>Parametric Pearson residuals test:</b> Assume Pearson residuals are Chi-squared distributed. For complex models, difficult to define parametric null distribution (unclear residual degrees of freedom).  <b>Nonparametric Pearson residuals test:</b> Parametric bootstrapping of the model to generate a nonparametric estimate of the null distribution of the Pearson statistic. Computational costly.	msme::P__disp()	GLMs	Hilbe & Robinson (2025)
			DHARMA::testDispersion(..., type="Pearson")	GLMs/GLMMs (naïve residual <i>df</i> )	Hartig (2024)
			performance::check_overdispersion()	GLMs/GLMMs (naïve residual <i>df</i> )	Lüdecke (2021)
			RVAideMemoire::overdisp.glmmer()	GLMMs (from lme4 package, naïve residual <i>df</i> , calculates only dispersion statistic, no test)	Herve (2025)
			DHARMA::testDispersion(..., refit=T, type="Pearson")	GLMs/GLMMs	Hartig (2024)
<b>Response variance</b>	Compares the expected to the observed variance in the response variable.	Expected variance of response variable calculated through simulations of fitted model. Fast nonparametrics but possibly less exact than working on the residual dispersion.	DHARMA::testDispersion(..., type="DHARMA")	GLMs/GLMMs	Hartig (2024)

## 154 **A short review of existing approaches to dispersion tests**

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

### 163 *Likelihood ratio tests*

164           A first general strategy for detecting dispersion problems is to compare a model  
165 with fixed dispersion to its nearest “relative” with variable dispersion using a likelihood  
166 ratio test (LRT) or another model selection technique, such as AIC (Yang et al., 2007).  
167 For count data, a practical example would be to compare a Poisson GLM as a null  
168 hypothesis to a negative binomial or generalised Poisson GLM (J. M. Hilbe, 2014), or  
169 to compare a binomial GLM to a beta-binomial GLM (Dunn & Smyth, 2018). While  
170 relatively easy to implement, the downside of this approach, apart from the higher  
171 computational costs resulting from fitting two models, is that it doesn’t provide any  
172 direct diagnostics of over- or underdispersion, but only compares a base model against  
173 an alternative. The alternative model, however, might also fit better or worse for reasons  
174 other than a dispersion problem. Moreover, using LRTs for detecting dispersion  
175 problems has also been discouraged as it may provide unreliable results (Dean, 1992)  
176 because it tends to underestimate the evidence against the base model (Lawless, 1987).

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

### 179 *Score tests*

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

194 An equivalent test related to the score test under certain conditions is the  
195 regression-based overdispersion test proposed by Cameron & Trivedi (1990). Under a  
196 Poisson model, the squared deviation of the observations from their fitted mean, after  
197 subtracting the observation itself and scaling by the fitted mean, has expectation zero. In  
198 contrast, under the negative binomial, it increases systematically with the mean. This  
199 motivates an auxiliary regression of the transformed variable against the fitted mean,  
200 with a significant slope indicating extra-Poisson variation. The main advantage of this

201 test against other score tests is its ease of implementation: it can be carried out after  
202 fitting a standard Poisson GLM. However, similar to an LRT, the linear regression  
203 imposes a particular form of overdispersion as an alternative hypothesis, and therefore,  
204 seems less general than the test based on Pearson residuals described below.

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

### 208 *Tests based on residual dispersion*

209 A third class of testing approaches, arguably the most intuitive, directly  
210 calculates a test statistic or goodness-of-fit metric on standardised model residuals. The  
211 most widely used test of this kind is based on the sum of the model's Pearson residuals.  
212 As Pearson residuals divide the raw residuals by the expected residual standard  
213 deviation, a correctly specified model is expected to have a Pearson residual of around 1  
214 for each observation. A dispersion statistic is then defined as the sum of squared  
215 Pearson residuals divided by the residual degrees of freedom. Models with a so-defined  
216 dispersion statistic  $> 1$  are considered overdispersed, while dispersion statistics  $< 1$  are  
217 underdispersed. Sometimes, a modification of this metric is often recommended by  
218 replacing the sum of squared Pearson residuals with the model deviance, which is  
219 typically more readily available. However, as Venables & Ripley (2002) discuss, this  
220 metric should be avoided, as it often deviates from 1, even for correctly specified  
221 GLMs.

222 Defining dispersion via the Pearson statistic has the added advantage that for a  
223 GLM, the expected distribution under the null hypothesis of a correctly specified model  
224 asymptotically follows a Chi-square distribution (McCullagh, 1985). This allows a

225 straightforward construction of a hypothesis test, where we compare the Pearson  
226 statistic to the chi-squared distribution with the respective residual degrees of freedom  
227 (*df*). This test is referred to with different terminologies, such as Pearson chi-squared  
228 dispersion test, Pearson residuals-based test for overdispersion, or simply Pearson  
229 dispersion test. Hereafter, we refer to this test as the **parametric Pearson residuals**  
230 **test**, to differentiate it from the nonparametric test based on Pearson residuals, discussed  
231 below.

232         An alternative approach to constructing a dispersion test based on the Pearson  
233 dispersion statistic involves generating a null distribution through parametric  
234 bootstrapping. A parametric bootstrap means that new data is simulated from the fitted  
235 model, and then the statistic of interest (in this case: the Pearson statistic of a fitted  
236 model) is calculated based on this data. The parametric bootstrap has been previously  
237 used for hypothesis tests in mixed-effects models where parametric null distributions  
238 were difficult to obtain (e.g., Barr et al., 2013; Luke, 2017), and thus it seems a logical  
239 alternative for more complicated models where the Chi-square distribution of the  
240 Pearson dispersion statistic cannot be taken for granted (see methods for GLMMs  
241 below). Nevertheless, implementing parametric bootstrapping in complex models can  
242 be less efficient for at least two reasons: it is time-consuming and prone to errors in  
243 model refits (Luke, 2017; Moral et al., 2017). A dispersion test based on this principle  
244 was implemented in R by Hartig (2024). Hereafter, we will refer to this test as the  
245 **nonparametric Pearson residuals test**.

#### 246 *Tests based on response variable variance*

247         Simulation approaches can also be useful to generate null distributions for  
248 alternative metrics of dispersion. A last class of dispersion test approaches, which, to

249 our knowledge, was introduced in the DHARMA R package (Hartig, 2024) but has not  
250 been discussed in the literature so far, involves defining a test statistic based on the  
251 dispersion of the response variable, rather than the residuals. More specifically, the test  
252 compares the observed data variance with the simulated data variances (which can be  
253 created conditional or unconditional on the fitted random effects for GLMMs). The  
254 dispersion statistic is then defined as the ratio between the observed variance and the  
255 mean simulated variance. Similar to the Pearson statistic, a ratio  $> 1$  indicates  
256 overdispersion, a ratio  $< 1$  indicates underdispersion, and a significance test is  
257 constructed based on the distribution of simulated variances.

258         From a theoretical viewpoint, this approach seems less elegant compared to the  
259 idea of using Pearson residuals, because the latter, by “standardising” the residual  
260 dispersion with the expected dispersion, allows each data point to contribute similarly to  
261 the dispersion statistic. In contrast, the test on the unstandardized response variable will  
262 be more influenced by large data points. However, the primary advantage of this  
263 approach is computational, as it enables the creation of a nonparametric estimate of the  
264 test statistic without requiring a re-fit of the model (in contrast to the nonparametric  
265 Pearson residuals test). Hereafter, we will refer to this test as the **simulation-based**  
266 **response variance test** to differentiate it from the tests based on Pearson residuals.

## 267 **Methods**

### 268 *Selected models and setup of the performance comparisons*

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

- 272 (1) The parametric Pearson residuals test  
273 (2) The nonparametric Pearson residuals test  
274 (3) The simulation-based response variance test

275 To compare the performance of these three tests, we simulated datasets based on  
276 the two main distributions that often present over- or underdispersion problems: the  
277 Poisson and the binomial (N/K) proportions. We varied the sample size (from 10 to  
278 10,000, depending on the simulation) and intercept (from -3 to 3, at the link function  
279 scale) of the simulated data for both distributions. We simulated a gradient of  
280 overdispersed data by adding noise to the linear predictor with values from a Gaussian  
281 distribution with a mean of zero and ten standard deviation values varying from 0 to 1.  
282 We evaluated the performance of the tests by comparing type I error, power, and  
283 dispersion statistics for all combinations of parameters in the simulated datasets.

284 All models were fitted using the functions *glm* from the stats package or *glmer*  
285 from the lme4 package (Bates et al., 2015) in R (v4.4; R Core Team, 2024). All  
286 dispersion tests were performed with the DHARMA package (Hartig, 2024). For the  
287 simulation-based response variance test and the nonparametric Pearson residuals tests,  
288 we set the number of simulations fixed at 250 (the default parameter in DHARMA). All  
289 simulations and analysis codes are available at this repository  
290 ([https://anonymous.4open.science/r/dispersion\\_test\\_GLMM/README.md](https://anonymous.4open.science/r/dispersion_test_GLMM/README.md)). The  
291 supplementary material provides a script file with instructions and examples for  
292 applying dispersion tests using the DHARMA package.

### 293 *Theoretical expectations*

294 The classical (1) parametric Pearson residuals test assumes that the sample size  
295 (n-asymptotic) and the expected values are sufficiently large (phi-asymptotic) (Venables

296 & Ripley, 2002). This implies that, when the expected counts (or intercept) and/or the  
297 number of observations are small, Pearson residuals may not provide reliable  
298 information about model fit (see S1). Some corrections for Pearson residuals with small  
299 sample sizes were suggested (e.g., Cordeiro, 2004; Cordeiro & Simas, 2009), but they  
300 are not currently implemented in the most common packages in R. Therefore, we expect  
301 the parametric Pearson residuals test to perform well for GLMs, except for very small  
302 sample sizes and expected counts (hereafter “small-data” situations).

303         Moreover, it is unclear whether the parametric Pearson residuals test approach  
304 can be extended to GLMMs or other hierarchical models, where counting the residual  
305 degrees of freedom (*df*) is not straightforward (Bolker et al., 2009; Luke, 2017). In  
306 mixed-effects models, the *df* used by a random effect are data-specific (adaptive  
307 shrinkage) and expected to be somewhere between one and the number of grouping  
308 factors (Baayen et al., 2008; Bolker et al., 2009; Luke, 2017). There exist approaches to  
309 approximate *df* for random effects in LMMs (e.g. Schaalje et al., 2002), but their  
310 generalisation to GLMMs is still an area of active research. Current R packages that  
311 implement the parametric Pearson residuals test approximate the *df* by the so-called  
312 naïve *df* (e.g.,  $n = 1$  per random effect) for testing LMMs/GLMMs (Table 1). We expect  
313 that the error imposed by this approximation increases with the number of random  
314 effect groups. To test this, we varied the number of groups in the random intercept (10,  
315 50, and 100 groups) of our simulated data.

316         In contrast to the parametric Pearson residual test, we expect the (2)  
317 nonparametric Pearson residuals test to be robust to small-data problems as well as the  
318 presence of random effects, as it doesn't rely on a particular parametric distribution.  
319 However, since the test uses parametric bootstrapping, we expected it to run much

320 slower than the other tests, especially for more complex GLMMs. For this purpose, we  
321 compared the runtime of the tests with a small set of simulated data (see S6).

322 For GLMMs tested with the (3) simulation-based response variance test, we  
323 compared the performance of the test under the two simulation approaches, conditional  
324 and unconditional to random effects. We expect to see lower power for the  
325 unconditional simulation results, as overdispersion is a phenomenon at the level of the  
326 model distribution (i.e., at a higher level). We evaluated the circumstances under which  
327 this test is reliable as a fast alternative to both dispersion tests based on Pearson  
328 residuals.

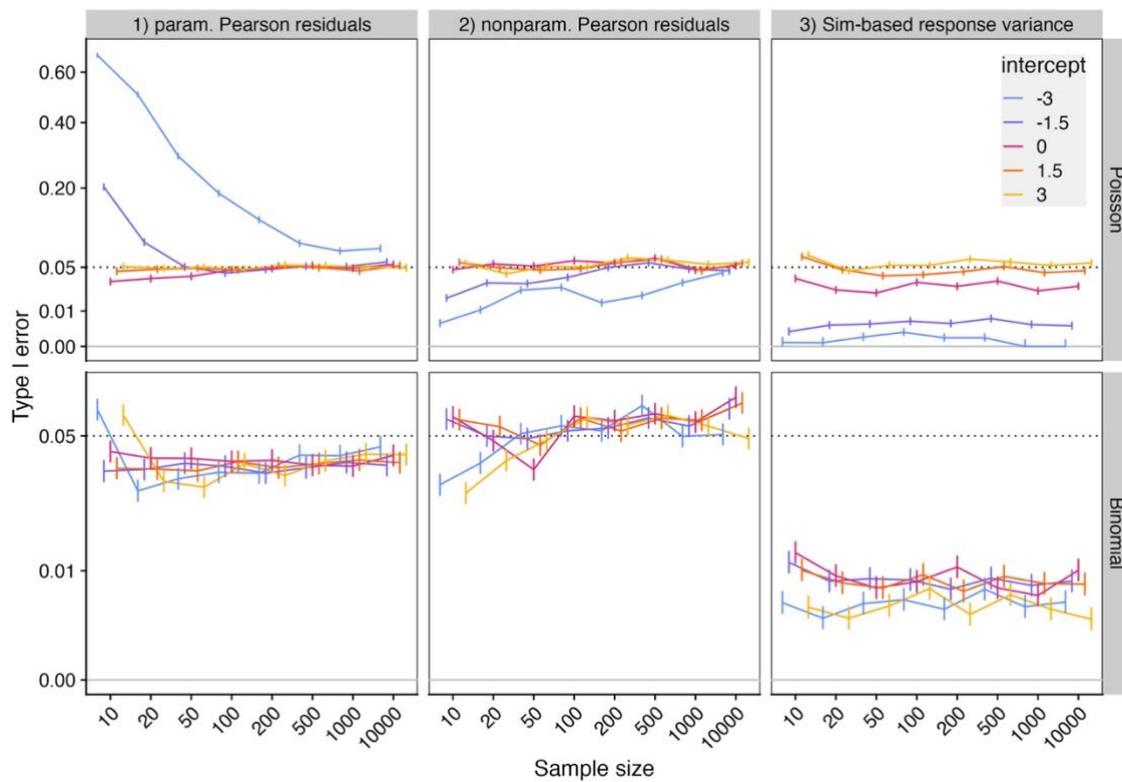
## 329 **Results**

### 330 *Performance on Poisson and binomial GLMs*

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

340 For binomial GLMs, the type I error rates for the parametric Pearson residuals  
341 test were generally conservatively calibrated around 0.04 (Fig. 2). Type I error rates for  
342 the nonparametric Pearson residuals test averaged around 0.05 and 0.06, except for the

343 very low and very high intercepts (-3 and 3). For the simulation-based response  
 344 variance test, type I error rates were conservatively very low for all simulated  
 345 parameters, bouncing below 0.01.



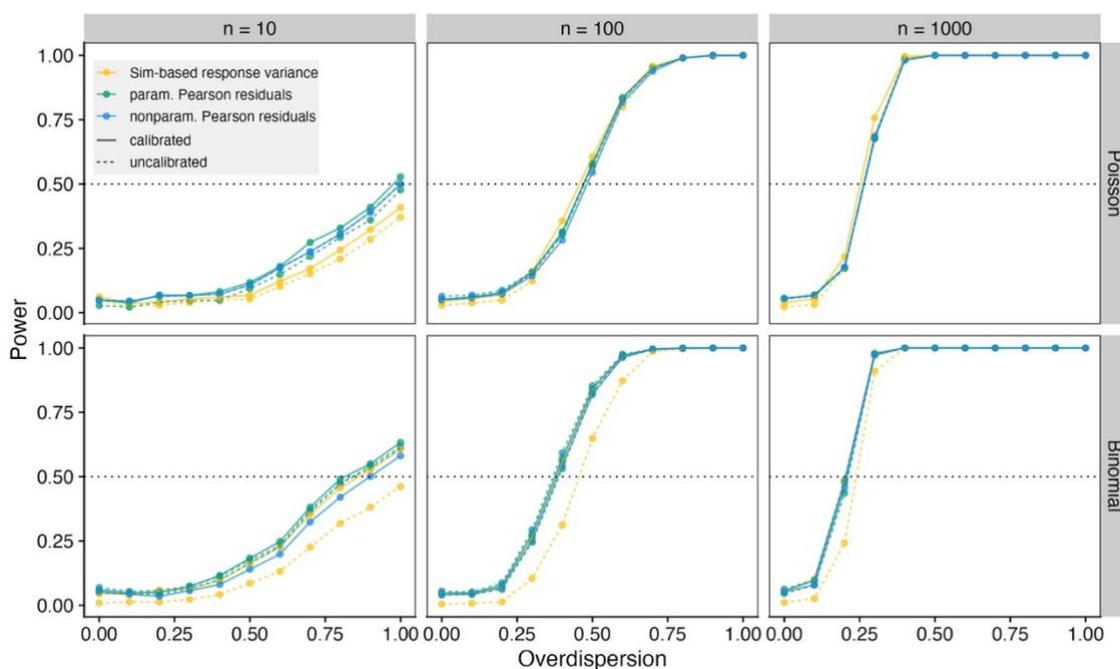
346

347 **Figure 2.** Simulation-based response variance tests have more conservative type I error  
 348 rates than both Pearson residuals tests. The three dispersion tests were applied to  
 349 Poisson (upper panels) and binomial proportion (lower panels) GLMs: 1a) parametric  
 350 Pearson residuals test, 1b) nonparametric Pearson residuals test, and 2) simulation-  
 351 based response variance test (see Table 1 for explanations). Simulations under different  
 352 sample sizes (x-axis) and intercepts (colours, values at the link function scale). In B),  
 353 the model is a binomial proportion with ten trials. All points include a 95% confidence  
 354 interval calculated based on exact binomial tests for the 10,000 simulations. Note the  
 355 square-root scale of the y-axis in plot A. The Dotted horizontal black line shows the  
 356 0.05 nominal value for type I error.

357 The statistical power of the simulation-based response variance test was lower  
 358 than the parametric and nonparametric Pearson residuals tests for both binomial and  
 359 Poisson GLMs, but tended to be similar with larger sample sizes (Fig. 3). We found that  
 360 the reason for this is the very conservative type I error rates (Fig. 2). When power is

361 calibrated by using the p-value at the 5% quantile of its empirical distribution for each  
 362 simulation (details in S5), the differences disappear (Fig. 3).

363 The dispersion statistics of the simulation-based response variance test were  
 364 highly dependent on the intercept, slope, and number of trials for the binomial model  
 365 (see S4, Fig. S4.2), and they tended to be smaller than those based on the Pearson  
 366 residuals. In contrast, for Poisson models, the values tended to be larger than those of  
 367 Pearson statistics (Fig. S3.5). This may also explain the lower uncorrected power for the  
 368 simulation-based response variance test, especially for binomial models.



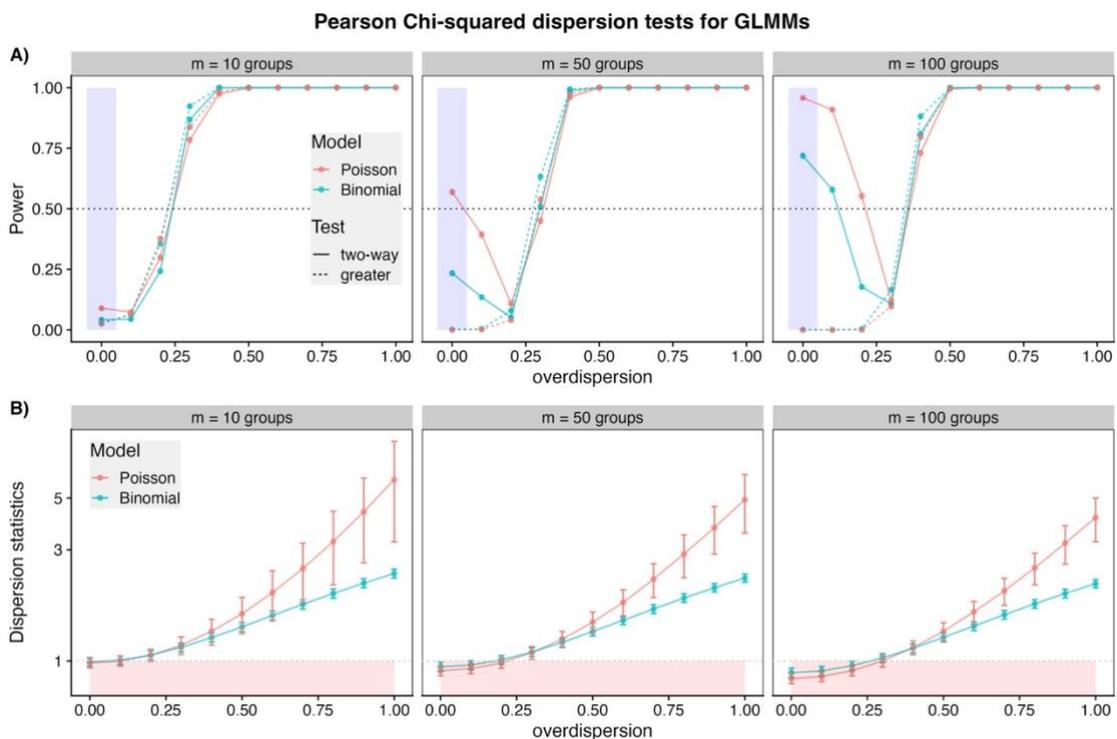
369

370 **Figure 3.** The simulation-based response variance test (in yellow) has lower power than  
 371 both Pearson residuals tests (green and blue) for GLMs unless power is calibrated by  
 372 type I error rates (dashed lines). Lower power is more evident for binomial models  
 373 (upper panel) and smaller sample sizes (first two columns). Results based on 10,000  
 374 simulations per combination of parameters for an intercept = 0 and slope = 1. For all  
 375 simulation results, see Fig. S5.1 and S5.2.

376 *GLMM performance*

377 For the GLMMs, we first compared the performance of the parametric Pearson  
 378 residuals test (two-sided) for an increasing number of groups ( $m$ ) in the random

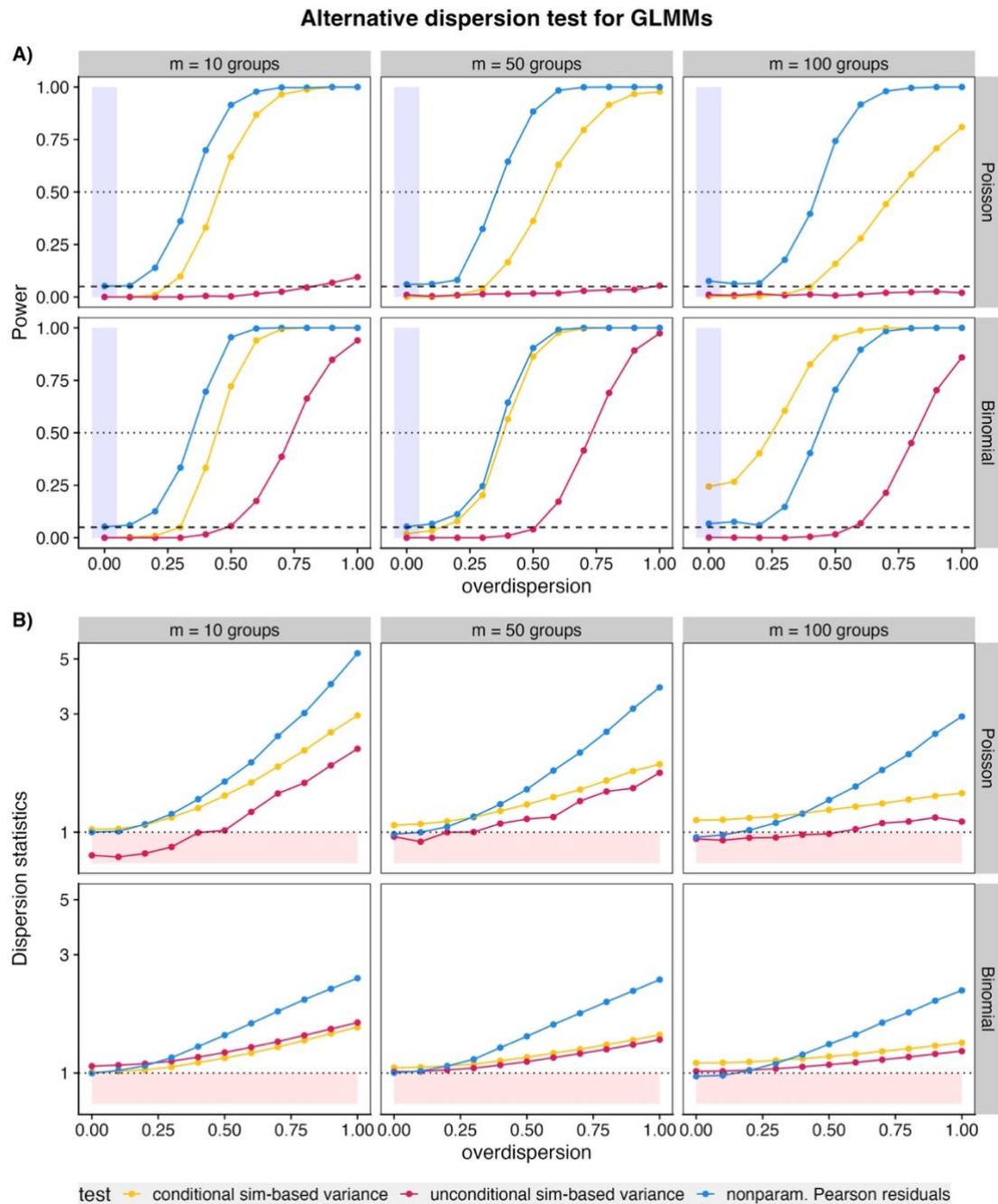
379 intercepts. As expected, the performance of the test failed for a large number of groups  
 380 in the random effects (Fig. 4A). The dispersion statistic was underestimated, and the  
 381 type I error rates were too high because the test detected significant underdispersion.  
 382 Testing only for overdispersion (“greater” test) when using the parametric Pearson  
 383 residuals test appears to be the only reasonable approach for GLMMs (Fig. 4A). Still, it  
 384 doesn’t prevent the dispersion statistics from being biased to lower values.



385  
 386 **Figure 4.** The parametric Pearson residuals test failed for GLMMs with many groups in  
 387 the random intercepts (plot panels). A) Power and type I error rates (blue shaded area)  
 388 for the “two-sided” (solid lines) and “greater” (dotted lines) Chi-squared tests for the  
 389 Pearson statistic. B) Pearson dispersion statistics with the red shaded area indicating  
 390 dispersion statistics estimated below 1 (underdispersion). Notice that the y-axis of plot  
 391 B is on a logarithmic scale of 10. Results with 10,000 simulations for an intercept of 0  
 392 and a sample size (n) of 1,000 data points.

393 When comparing the alternative dispersion tests for GLMMs, the nonparametric  
 394 Pearson residuals test presented very good results, with a type I error rate around 0.05  
 395 (Fig. S6.1 and S6.2) and higher power than the simulation-based response variance tests  
 396 (Fig. 5). As expected, the unconditional simulation-based response variance test had the

397 worst performance: very low type I errors (Fig. S6.1 and S6.2), very low power, and  
398 dispersion statistics below 1 (Fig. 5B), especially for Poisson models. The conditional  
399 simulation-based response variance test also had very small type I errors (Fig. S6.1 and  
400 S6.2), but power increased with the simulated overdispersion. The performance of both  
401 simulation-based response variance tests (unconditional and conditional) didn't change  
402 much with the number of groups for the Poisson GLMMs, but it improved for the  
403 binomial GLMMs with the increasing number of groups in the random intercept.



404

405 **Figure 5.** The nonparametric Pearson residuals test showed correct Type 1 error, higher  
 406 power, and larger dispersion statistics than the simulated-based response variance tests  
 407 (conditional and unconditional to all random effects) for Poisson and binomial GLMMs.  
 408 Power (A), type I error (shaded blue area in A), and dispersion statistics (B) for the  
 409 alternative dispersion tests for Poisson and binomial GLMMs with different numbers of  
 410 groups in random intercepts. The dashed horizontal line in (A) indicates the nominal  
 411 value of 0.05 for type I error. The dotted horizontal line in (A) indicates the 50% power,  
 412 and the dotted horizontal line in (B) indicates the dispersion statistics of 1. The results  
 413 are based on 1,000 simulations per combination of parameters, with an intercept of 0  
 414 and a sample size ( $n$ ) of 1,000.

415 **Discussion**

416           The goal of this study was to find a dispersion test that is widely applicable  
417 across different GLM and GLMM distributions and random-effects structures. Our  
418 conclusion is that the nonparametric Pearson residuals test is the most reliable general  
419 test currently available. For GLMs, this test exhibited similar power as the parametric  
420 Pearson residuals test but with more reliable type I error rates in small-data situations.  
421 The downside of this test is that it can be computationally expensive, with runtimes in  
422 the order of minutes for larger GLMMs.

423           The two alternative tests that we considered have advantages in particular  
424 situations. The simulation-based response variance test for GLMs is fast to compute, but  
425 has a dispersion statistic that is more difficult to interpret and often too conservative  
426 type I errors. This resulted in low power if not additionally calibrated by a simulated p-  
427 value distribution. The parametric Pearson residuals test is computationally efficient,  
428 but it is unreliable in small-data situations and in the presence of random effects. Below,  
429 we discuss these points in more detail and provide recommendations for general users  
430 who rely on already implemented R packages for model fit and diagnostics.

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

432           We showed that the parametric Pearson residuals test, although popular, quick,  
433 and relatively easy to compute, has two main disadvantages: it does not perform well in  
434 (1) small-data situations (Fig. 2) and (2) in the presence of random effects (Fig. 4). The  
435 reason for the first problem can be attributed to the mismatch between the Pearson  
436 statistic distribution and the Chi-squared distribution under small-data conditions (Fig.  
437 S1.1 and S1.2). This phenomenon has already been studied (e.g., Fletcher, 2012; Kuss,  
438 2002), with suggested corrections (Farrington, 1996; McCullagh, 1985). However, none  
439 of these corrections are implemented in the current R packages (Table 1), and we

440 believe that it will be difficult to devise corrections that work across a wide range of  
441 distributions.

442         The reason for the second problem is that counting 1 degree of freedom (*df*) for  
443 a random effect, as done in most implementations of this test, is typically an  
444 underestimation of the true model *df*, which increases in magnitude with the increasing  
445 number of levels of the random effect. The result is a bias in the dispersion statistic  
446 towards underdispersion that increases with the number of random effect levels (Fig. 4).  
447 Two-sided tests would therefore often wrongly detect significant underdispersion  
448 problems in perfectly valid GLMMs, which is likely the reason why most R  
449 implementations of this test only test for overdispersion. When applying this test for  
450 GLMMs, we recommend following the same approach and ignoring dispersion statistics  
451 smaller than 1. Nevertheless, it is an unsatisfactory solution since the biased dispersion  
452 statistic will also cause a loss of power.

453         A possible solution for GLMMs could be using a better approximation of the  
454 residual degrees of freedom (*df*). For LMMs, approximations for denominator *df* have  
455 been successfully used for hypothesis testing (Luke, 2017), for example, the  
456 Satterthwaite (1946) and the Kenward-Roger (2009). Although there is some evidence  
457 that these approximations are also accurate for GLMMs (Stroup, 2015), the main R  
458 packages implementing some of these approximation methods are currently limited to  
459 LMMs (e.g., *pbkrtest* Halekoh & Højsgaard, 2014; *lmerTest* Kuznetsova et al., 2017).  
460 However, the recently released package *glmmrBase* (Watson, 2024) allows these  
461 methods to be applied to GLMMs. We performed some parametric Pearson residuals  
462 tests for Poisson GLMMs using a modified residual *df* approximation (see S8).  
463 Although the parametric Pearson residuals tests with the approximated residual *df*  
464 performed much better than those with the naïve residual *df*, they still underperformed

465 compared to the nonparametric Pearson test when having a large number of groups in  
466 the random effects (Fig. S8.4), especially for small-data situations.

467 *When are simulation-based response variance tests an alternative?*

468         The simulation-based response variance test developed in the R package  
469 DHARMA (Hartig, 2024) is the main alternative to the family of Pearson residuals tests.  
470 Its principle is simple: when the model is correctly specified, the variance of the  
471 observed data should match the variance of the data simulated from the model. The  
472 main advantage of this approach is that it is a non-parametric test that can be applied to  
473 any model structure and it does not require refitting the model, which makes it both  
474 considerably faster and easier to implement in statistical software. We also note that for  
475 GLMMs, simulations should be performed conditionally to avoid a loss of power,  
476 presumably due to the increased variability created by re-simulating the random effects  
477 (unconditional simulations).

478         The disadvantages of this approach are that it is often overly conservative,  
479 resulting in lower power compared to the Pearson residuals tests. Additionally, the  
480 calculated dispersion statistic differs from the Pearson dispersion statistic, making it  
481 difficult to compare the two approaches. We conjectured that both problems could be  
482 related to the fact that the test statistic is based on the raw variance (and not a scaled  
483 variance, as for the Pearson statistics), and therefore observations with large values may  
484 be overrepresented in the statistics. We considered scaling each observation with  
485 expected variance, but this is not readily available for a wide class of models, and using  
486 simulations to approximate it fails for discrete-valued distributions (see S7).

487 **Conclusions and recommendations**

488 In conclusion, while neither of the considered options excelled in all dimensions  
 489 (Fig. 6), our base recommendation is that for standard GLMs with sufficient data, the  
 490 parametric Pearson Chi-square test, available in many packages, can be safely used. In  
 491 complex situations, particularly for GLMMs, we recommend the nonparametric Pearson  
 492 residuals test. It has very few weaknesses, other than being computationally costly to  
 493 calculate. If the nonparametric Pearson residuals test cannot be calculated due to speed  
 494 or convergence problems with refitting complex models, we recommend using the  
 495 simulation-based response variance test with simulations performed conditionally on the  
 496 fitted random effects. All three approaches are available via the *testDispersion* function  
 497 in the DHARMA R package (Hartig, 2024). We provide a tutorial with instructions and  
 498 an example for applying dispersion tests using the DHARMA package on the repository  
 499 website (<https://theoreticalecology.github.io/dispersionTest/>).

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

500

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

507           Although our simulation examples concentrated on overdispersion, the tests  
508 under consideration in our study can equally be used to detect underdispersion problems  
509 by testing the dispersion “two-sided” or “less than” against null statistics. The clear  
510 exception would be testing for underdispersion using the parametric Pearson residuals  
511 test for GLMMs, which would be anti-conservative due to the discussed bias towards  
512 underdispersion in the presence of random effects.

513 *Recommendations for practical data analysis when using dispersion tests*

514           For the interpretation and applied data analysis, we stress that getting a  
515 significant over-/underdispersion result does not necessarily indicate that the  
516 distribution must be changed. First, hypothesis tests famously evaluate statistical rather  
517 than practical significance. In other words, a significant test for overdispersion indicates  
518 that the overdispersion signal deviates from a null expectation, but the p-value does not  
519 measure the strength of the deviation. The first step in a dispersion test should thus be to  
520 examine how much the dispersion statistic deviates from the expected value of 1. For  
521 very large sample sizes, small departures from 1 may be statistically significant, but  
522 they may not necessarily warrant a change to the model. Second, after finding that a  
523 dispersion problem is both significant and meaningful, we suggest first checking for  
524 problems other than the distribution, such as heteroscedasticity, missing predictors,  
525 incorrect link function, excess of zeros, or overfitting. In our experience, these types of  
526 model misspecifications often cause over-/underdispersion, but can be distinguished  
527 from a “real” distributional problem through careful residual checks. Blindly changing  
528 the distribution only masks the problem, without offering a real remedy to the  
529 underlying problems.

530 Finally, after having convinced ourselves through these previous investigations  
531 that we are facing an ‘intrinsic’ under-/overdispersion problem, we should consider  
532 changing the GLM distribution. A traditional and flexible solution is using the ‘quasi’  
533 distributions (Wedderburn, 1974), which essentially correct p-values, but have the  
534 disadvantage that they do not represent an explicit data-generating process with  
535 associated likelihood, which does not allow, for example, to simulate from the fitted  
536 model. A second alternative to add dispersion is using observation-level random effects  
537 (Bolker et al., 2009; Elston et al., 2001; Harrison, 2014; Ozgul et al., 2009). While often  
538 offering a reasonable solution, we feel that the excessive use of REs tends to create  
539 problems in the calculation of other statistical indicators (such as p-values) that we  
540 would rather avoid. For that reason, we feel the best solution to address ‘intrinsic’  
541 under-/overdispersion is to switch to the corresponding variable-dispersion  
542 distributions, such as the negative binomial (Harrison, 2014) for overdispersed and the  
543 Conway-Maxwell-Poisson distribution (Lynch et al., 2014) for underdispersed Poisson  
544 models, or the beta-binomial distribution for overdispersed binomial models (Harrison,  
545 2015). Regardless of the approach, an “over-/underdispersion-free” GLM/GLMM is  
546 essential for better interpreting model results and facilitating sound scientific  
547 discoveries.

## 548 **References**

- 549 Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with  
550 crossed random effects for subjects and items. *Journal of Memory and*  
551 *Language*, 59(4), 390–412. <https://doi.org/10.1016/j.jml.2007.12.005>
- 552 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for  
553 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and*  
554 *Language*, 68(3), 255–278. <https://doi.org/10.1016/j.jml.2012.11.001>
- 555 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects  
556 Models Using lme4. *Journal of Statistical Software*, 67(1), 1–48.  
557 <https://doi.org/10.18637/jss.v067.i01>

- 558 Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H.  
559 H., & White, J.-S. S. (2009). Generalized linear mixed models: A practical guide  
560 for ecology and evolution. *Trends in Ecology & Evolution*, 24(3), 127–135.  
561 <https://doi.org/10.1016/j.tree.2008.10.008>
- 562 Cameron, A. C., & Trivedi, P. (2023). *overdisp: Overdispersion in count data multiple*  
563 *regression analysis* (Version 0.1.2) [Computer software]. [https://CRAN.R-](https://CRAN.R-project.org/package=overdisp)  
564 [project.org/package=overdisp](https://CRAN.R-project.org/package=overdisp)
- 565 Cameron, A. C., & Trivedi, P. K. (1990). Regression-based tests for overdispersion in  
566 the Poisson model. *Journal of Econometrics*, 46(3), 347–364.  
567 [https://doi.org/10.1016/0304-4076\(90\)90014-K](https://doi.org/10.1016/0304-4076(90)90014-K)
- 568 Campbell, H. (2021). The consequences of checking for zero-inflation and  
569 overdispersion in the analysis of count data. *Methods in Ecology and Evolution*,  
570 12(4), 665–680. <https://doi.org/10.1111/2041-210X.13559>
- 571 Collings, B. J., & Margolin, B. H. (1985). Testing Goodness of Fit for the Poisson  
572 Assumption When Observations Are Not Identically Distributed. *Journal of the*  
573 *American Statistical Association*, 80(390), 411–418.
- 574 Cordeiro, G. M. (2004). On Pearson's residuals in generalized linear models. *Statistics*  
575 *& Probability Letters*, 66(3), 213–219. <https://doi.org/10.1016/j.spl.2003.09.004>
- 576 Cordeiro, G. M., & Simas, A. B. (2009). The distribution of Pearson residuals in  
577 generalized linear models. *Computational Statistics & Data Analysis*, 53(9),  
578 3397–3411. <https://doi.org/10.1016/j.csda.2009.02.025>
- 579 Dean. (1992). Testing for Overdispersion in Poisson and Binomial Regression Models.  
580 *Journal of the American Statistical Association*, 87(418), 451–457.  
581 <https://doi.org/10.2307/2290276>
- 582 Dean, C., & Lawless, J. F. (1989). Tests for Detecting Overdispersion in Poisson  
583 Regression Models. *Journal of the American Statistical Association*, 84(406),  
584 467–472. <https://doi.org/10.1080/01621459.1989.10478792>
- 585 Dunn, P. K., & Smyth, G. K. (2018). *Generalized Linear Models With Examples in R*.  
586 Springer New York. <https://doi.org/10.1007/978-1-4419-0118-7>
- 587 Elston, D. A., Moss, R., Boulinier, T., Arrowsmith, C., & Lambin, X. (2001). Analysis  
588 of aggregation, a worked example: Numbers of ticks on red grouse chicks.  
589 *Parasitology*, 122(05), 563–569.
- 590 Farrington, C. P. (1996). On Assessing Goodness of Fit of Generalized Linear Models to  
591 Sparse Data. *Journal of the Royal Statistical Society. Series B (Methodological)*,  
592 58(2), 349–360.
- 593 Feng, C., Li, L., & Sadeghpour, A. (2020). A comparison of residual diagnosis tools for  
594 diagnosing regression models for count data. *BMC Medical Research*  
595 *Methodology*, 20(1), 175. <https://doi.org/10.1186/s12874-020-01055-2>
- 596 Fisher, R. A. (1950). The Significance of Deviations from Expectation in a Poisson  
597 Series. *Biometrics*, 6(1), 17–24. <https://doi.org/10.2307/3001420>
- 598 Fletcher, D. J. (2012). Estimating overdispersion when fitting a generalized linear  
599 model to sparse data. *Biometrika*, 99(1), 230–237.  
600 <https://doi.org/10.1093/biomet/asr083>

- 601 Halekoh, U., & Højsgaard, S. (2014). A Kenward-Roger Approximation and Parametric  
602 Bootstrap Methods for Tests in Linear Mixed Models – The R Package pbkrtest.  
603 *Journal of Statistical Software*, 59, 1–32. <https://doi.org/10.18637/jss.v059.i09>
- 604 Harrison, X. A. (2014). Using observation-level random effects to model overdispersion  
605 in count data in ecology and evolution. *PeerJ*, 2, e616.  
606 <https://doi.org/10.7717/peerj.616>
- 607 Harrison, X. A. (2015). A comparison of observation-level random effect and Beta-  
608 Binomial models for modelling overdispersion in Binomial data in ecology &  
609 evolution. *PeerJ*, 3, e1114. <https://doi.org/10.7717/peerj.1114>
- 610 Hartig, F. (2024). *DHARMa: Residual Diagnostics for Hierarchical (Multi-Level /*  
611 *Mixed) Regression Models* (Version 0.4.7) [Computer software].  
612 <https://CRAN.R-project.org/package=DHARMa>
- 613 Herve, M. (2025). *RVAideMemoire: Testing and plotting procedures for biostatistics*  
614 (Version 0.9-83-11) [Computer software]. [https://CRAN.R-](https://CRAN.R-project.org/package=RVAideMemoire)  
615 [project.org/package=RVAideMemoire](https://CRAN.R-project.org/package=RVAideMemoire)
- 616 Hilbe, J. M. (2011). *Negative Binomial Regression*.
- 617 Hilbe, J. M. (2014). *Modeling count data*. Cambridge University Press.
- 618 Hilbe, J., & Robinson, A. (2025). *msme: Functions and datasets for “methods of*  
619 *statistical model estimation”* (Version 0.5.4) [Computer software].  
620 <https://CRAN.R-project.org/package=msme>
- 621 Jackman, S. (2024). *pscl: Classes and methods for R developed in the political science*  
622 *computational laboratory* (Version 1.5.9) [Computer software]. University of  
623 Sydney. <https://github.com/atahk/pscl/>
- 624 Kenward, M. G., & Roger, J. H. (2009). An improved approximation to the precision of  
625 fixed effects from restricted maximum likelihood. *Computational Statistics &*  
626 *Data Analysis*, 53(7), 2583–2595. <https://doi.org/10.1016/j.csda.2008.12.013>
- 627 Kleiber, C., & Zeileis, A. (2008). *Applied econometrics with R [R package AER version*  
628 *1.2-14]*. Springer-Verlag. <https://doi.org/10.1007/978-0-387-77318-6>
- 629 Kuss, O. (2002). Global goodness-of-fit tests in logistic regression with sparse data.  
630 *Statistics in Medicine*, 21(24), 3789–3801. <https://doi.org/10.1002/sim.1421>
- 631 Kuznetsova, A., Brockhoff, P., & Christensen, R. (2017). lmerTest Package: Tests in  
632 Linear Mixed Effects Models. *Journal of Statistical Software, Articles*, 82(13).  
633 <https://doi.org/10.18637/JSS.V082.I13>
- 634 Lai, J., Lortie, C. J., Muenchen, R. A., Yang, J., & Ma, K. (2019). Evaluating the  
635 popularity of R in ecology. *Ecosphere*, 10(1), e02567.  
636 <https://doi.org/10.1002/ecs2.2567>
- 637 Lawless, J. F. (1987). Negative binomial and mixed Poisson regression. *Canadian*  
638 *Journal of Statistics*, 15(3), 209–225. <https://doi.org/10.2307/3314912>
- 639 Lopez-Quílez, V. G.-R. J. F.-F. A. (2005). Detecting clusters of disease with R. *Journal*  
640 *of Geographical Systems*, 7(2), 189–206. [https://doi.org/10.1007/s10109-005-](https://doi.org/10.1007/s10109-005-0156-5)  
641 [0156-5](https://doi.org/10.1007/s10109-005-0156-5)
- 642 Lüdecke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021).  
643 performance: An R package for assessment, comparison and testing of statistical

- 644 models. *Journal of Open Source Software*, 6(60), 3139.  
645 <https://doi.org/10.21105/joss.03139>
- 646 Luke, S. G. (2017). Evaluating significance in linear mixed-effects models in R.  
647 *Behavior Research Methods*, 49(4), 1494–1502. [https://doi.org/10.3758/s13428-](https://doi.org/10.3758/s13428-016-0809-y)  
648 016-0809-y
- 649 Lynch, H. J., Thorson, J. T., & Shelton, A. O. (2014). Dealing with under- and over-  
650 dispersed count data in life history, spatial, and community ecology. *Ecology*,  
651 95(11), 3173–3180. <https://doi.org/10.1890/13-1912.1>
- 652 McCullagh, P. (1985). On the Asymptotic Distribution of Pearson's Statistic in Linear  
653 Exponential-Family Models. *International Statistical Review / Revue*  
654 *Internationale de Statistique*, 53(1), 61–67. <https://doi.org/10.2307/1402880>
- 655 McCullagh, P., & Nelder, J. (1989). Generalized linear models. *Journal of the Royal*  
656 *Statistical Society*, 135(3), 370–384.
- 657 Moral, R. A., Hinde, J., & Demétrio, C. G. B. (2017). Half-Normal Plots and  
658 Overdispersed Models in R: The hnp Package. *Journal of Statistical Software*,  
659 81, 1–23. <https://doi.org/10.18637/jss.v081.i10>
- 660 Ohara Hines, R. J. (1997). A comparison of tests for overdispersion in generalized linear  
661 models. *Journal of Statistical Computation and Simulation*, 58(4), 323–342.  
662 <https://doi.org/10.1080/00949659708811838>
- 663 Ozgul, A., Oli, M. K., Bolker, B. M., & Perez-Heydrich, C. (2009). Upper respiratory  
664 tract disease, force of infection, and effects on survival of gopher tortoises.  
665 *Ecological Applications*, 19(3), 786–798.
- 666 Papadakis, M., Tsagris, M., Fafalios, S., Dimitriadis, M., & Lasithiotakis, M. (2025).  
667 *Rfast2: A collection of efficient and extremely fast R functions II* (Version  
668 0.1.5.4) [Computer software]. <https://CRAN.R-project.org/package=Rfast2>
- 669 Quine, M. P., & Seneta, E. (1987). Bortkiewicz's Data and the Law of Small Numbers.  
670 *International Statistical Review / Revue Internationale de Statistique*, 55(2),  
671 173–181. <https://doi.org/10.2307/1403193>
- 672 R Core Team. (2024). *R: a language and environment for statistical computing* (Version  
673 v4.4.1) [Computer software]. R Foundation for Statistical Computing.  
674 <https://www.R-project.org/>
- 675 Rao, C. R. (1948). Large sample tests of statistical hypotheses concerning several  
676 parameters with applications to problems of estimation. *Mathematical*  
677 *Proceedings of the Cambridge Philosophical Society*, 44(1), 50–57.  
678 <https://doi.org/10.1017/S0305004100023987>
- 679 Rhodes, J. R. (2015). Mixture models for overdispersed data. In G. A. Fox, V. J. Sosa, &  
680 S. M. Negrete-Yankelevich, *Ecological Statistics: Contemporary theory and*  
681 *applications*. Oxford University Press.
- 682 Satterthwaite, F. E. (1946). An Approximate Distribution of Estimates of Variance  
683 Components. *Biometrics Bulletin*, 2(6), 110–114.  
684 <https://doi.org/10.2307/3002019>
- 685 Schaalje, G. B., McBride, J. B., & Fellingham, G. W. (2002). Adequacy of  
686 approximations to distributions of test statistics in complex mixed linear models.

- 687 *Journal of Agricultural, Biological, and Environmental Statistics*, 7(4), 512–  
688 524. <https://doi.org/10.1198/108571102726>
- 689 Stroup, W. W. (2015). Rethinking the Analysis of Non-Normal Data in Plant and Soil  
690 Science. *Agronomy Journal*, 107(2), 811–827.  
691 <https://doi.org/10.2134/agronj2013.0342>
- 692 Student. (1919). An explanation of deviations from poisson’s law in practice.  
693 *Biometrika*, 12, 211.
- 694 Touchon, J. C., & McCoy, M. W. (2016). The mismatch between current statistical  
695 practice and doctoral training in ecology. *Ecosphere*, 7(8), e01394.  
696 <https://doi.org/10.1002/ecs2.1394>
- 697 Venables, W. N., & Ripley, B. D. (2002). *Modern Applied Statistics with S*. Springer  
698 New York. <https://doi.org/10.1007/978-0-387-21706-2>
- 699 Watson, S. I. (2024). *Generalised Linear Mixed Model Specification, Analysis, Fitting,  
700 and Optimal Design in R with the glmmr Packages* (No. arXiv:2303.12657).  
701 arXiv. <https://doi.org/10.48550/arXiv.2303.12657>
- 702 Wedderburn, R. W. M. (1974). Quasi-likelihood functions, generalized linear models,  
703 and the Gauss—Newton method. *Biometrika*, 61(3), 439–447.  
704 <https://doi.org/10.1093/biomet/61.3.439>
- 705 Xekalaki, E. (2014). On the distribution theory of over-dispersion. *Journal of Statistical  
706 Distributions and Applications*, 1(1), 19. [https://doi.org/10.1186/s40488-014-  
0019-z](https://doi.org/10.1186/s40488-014-<br/>707 0019-z)
- 708 Yang, Z., Hardin, J. W., Addy, C. L., & Vuong, Q. H. (2007). Testing Approaches for  
709 Overdispersion in Poisson Regression versus the Generalized Poisson Model.  
710 *Biometrical Journal*, 49(4), 565–584. <https://doi.org/10.1002/bimj.200610340>
- 711 Zeileis, A., & Hothorn, T. (2002). Diagnostic checking in regression relationships. *R  
712 News*, 2(3), 7–10.
- 713

# Supplementary Material of

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

### methods comparison and practical guide

Melina de Souza Leite<sup>1\*</sup>, Daniel Rettelbach<sup>1,2</sup> & Florian Hartig<sup>1</sup>

1. Theoretical Ecology, University of Regensburg, Germany

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

\*corresponding author: [melina.souza-leite@ur.de](mailto:melina.souza-leite@ur.de)

#### S1. Pearson statistics and Chi-squared distribution

For GLMs, the parametric Pearson residuals test assumes that the sample size (n-asymptotic) and the expected values are sufficiently large (phi-asymptotic).

Therefore, when the expected counts (or intercept) and/or the number of observations are small, Pearson residuals may not provide reliable information about model fit. To

test boundaries where Pearson statistics fail, we simulated data with very different

sample sizes (from 10 to 10,000, depending on the simulation) and intercepts (from -3

to 3, at the link function scale) for Poisson and binomial proportion GLMs. For each

distribution and parameter combination, we used the Kolmogorov-Smirnov test (KS

test) of adherence to compare the empirical distribution of 1000 simulations of the

Pearson residuals with the Chi-squared distribution having the same residual degrees of

freedom. We repeated this procedure 100 times and recorded the proportion of

significant KS tests.

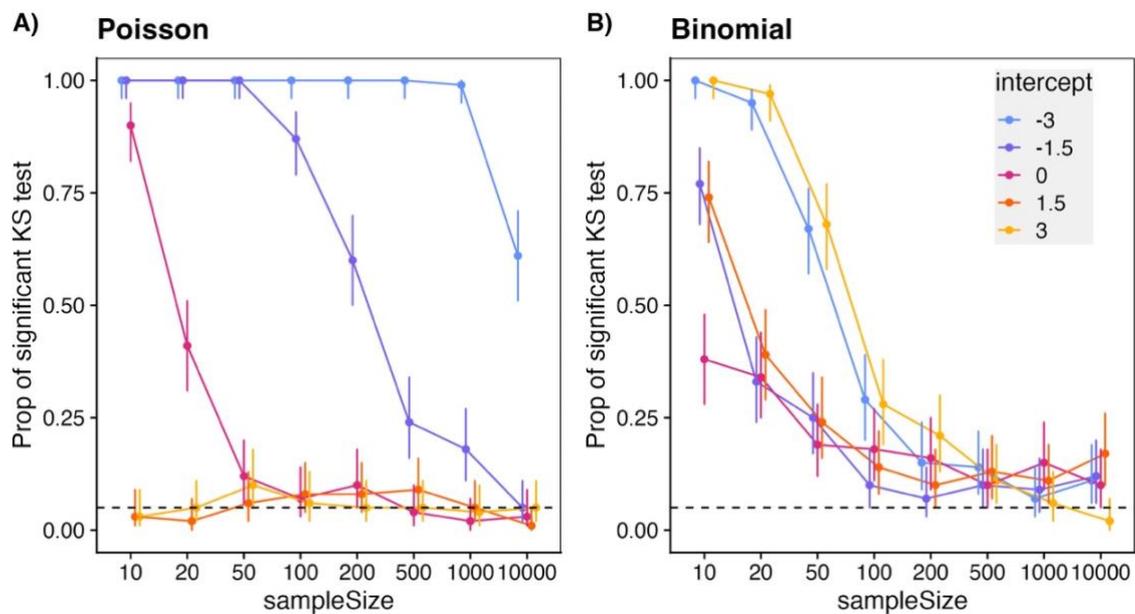
For the Poisson GLMs, the Pearson statistics distribution clearly departed from

the Chi-square distribution for very small intercepts (-3, -1.5) and sample sizes (10, 20

and 50) (Figure S1.1 A). Even for very large sample sizes (10,000), the distribution did

25 not approximate the Chi-squared distribution for the smallest simulated intercept (-3).  
 26 Consequently, the KS tests showed all significant results for all simulations with the  
 27 intercept at -3, except for the largest sample size (10,000), where it decreased to 60%.  
 28 As expected, the proportion of significant results decreased with sample size for  
 29 intercepts at -1.5 and 0. For larger intercepts, it remained around 5% for all sample sizes  
 30 (Figure S1.2A).

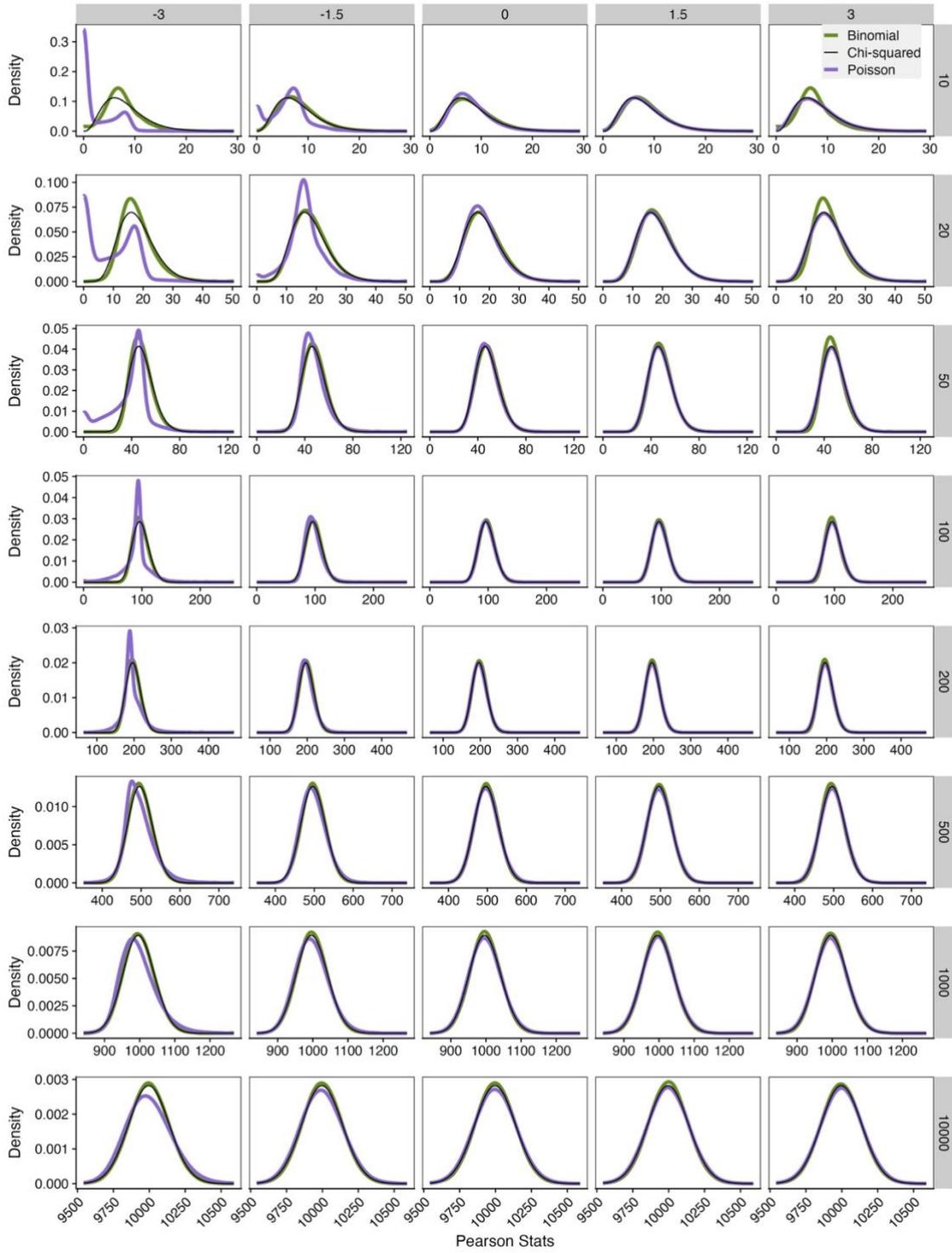
31 For the **binomial GLMs**, the Pearson statistics distribution clearly departed  
 32 from the Chi-squared distribution for very small and large intercepts (-3, 3) and small  
 33 sample sizes (10, 20, 50) (Figure S1.1B). The proportion of significant KS tests  
 34 decreased with sample size, but did not reach the nominal value of 0.05, even for very  
 35 large sample sizes and intermediate intercept values (-1.5, 0, 1.5).



36

37 **Figure S1.1.** Proportion of significant Kolmogorov-Smirnov adherence tests between  
 38 the empirical distribution of 1000 simulations of the Pearson statistics and a Chi-  
 39 squared distribution with the same residual degrees of freedom for A) Poisson and B)  
 40 binomial GLMs. Proportions were calculated from 100 simulations for each  
 41 combination of the data parameters (sample size and intercept). For binomial data, the  
 42 number of trials was fixed at 10. The 95% confidence intervals (vertical lines) were  
 43 drawn from binomial exact tests for each result with  $p = 0.05$ .

**Pearson Statistics X Chi-squared distribution**



44

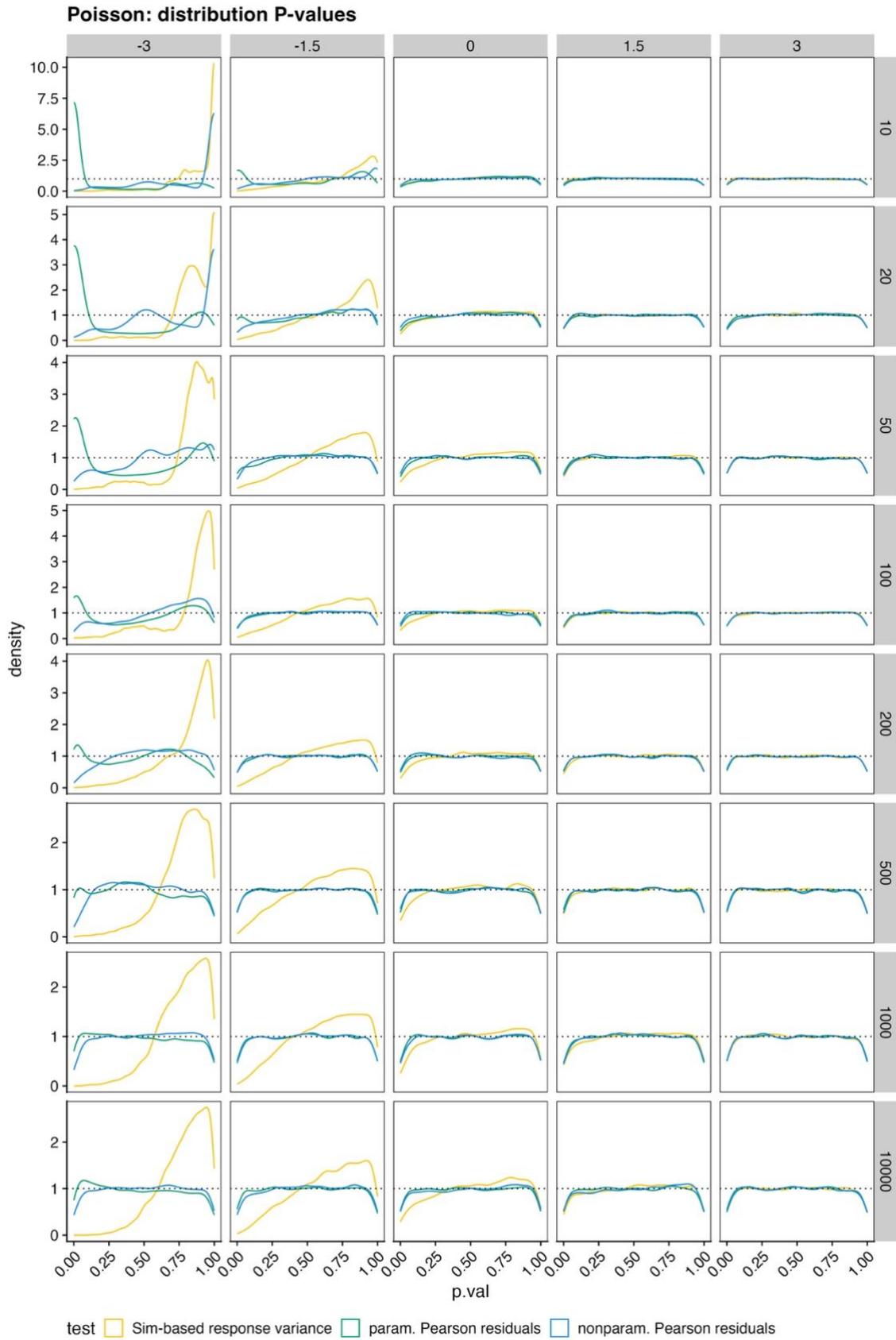
45 **Figure S1.2.** Mean Pearson statistics distribution (from 100 simulated curves) for the  
 46 binomial (green) and Poisson (purple), and the Chi-square distribution in black.

## 47 **S2. Type I error rates for the GLMs**

48           Figures S2.1 and S2.2 show the distribution of the p-values for the dispersion  
49 tests applied to the Poisson and binomial GLMs, respectively, with 10,000 simulations  
50 for each combination of intercept and sample size. For the dispersion tests with correct  
51 type I error rates around the nominal value of 0.05, the distributions of p-values should  
52 present a uniform distribution with density 1.

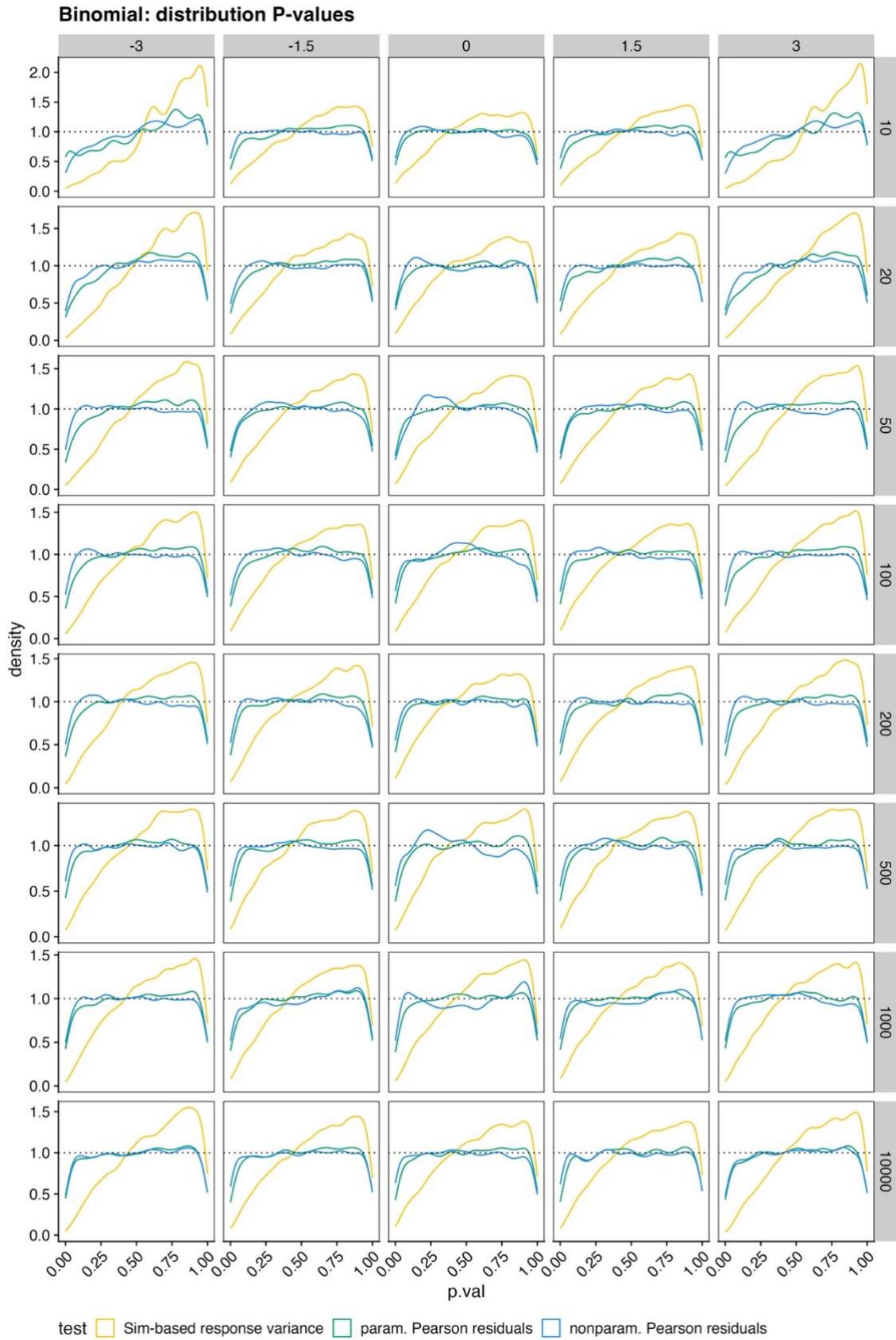
53           For the Poisson GLMs (Figure S2.1), the simulation-based response variance  
54 test (in red) presented the largest departure of the expected distribution for the smallest  
55 intercepts (-3, -1.5) across all sample sizes. This explains why the type I error rates for  
56 the simulation-based residual tests were so low and varied according to the intercept but  
57 didn't change with the sample size (main text Figure 2A). The parametric Pearson test  
58 had the opposite pattern with very low p-values for the smallest intercept (-3), but it  
59 tended to approximate the uniform distribution (decreasing the peak for the low p-  
60 values) with sample size. The p-values for the nonparametric Pearson test also showed a  
61 departure from the uniform distribution for the smallest intercept (-3), but tended to  
62 approach the uniform distribution with larger sample sizes and intercepts.

63           For the binomial GLMs (Figures S2.2), the p-values distribution of the  
64 simulation-based response variance test also presented the largest departure from the  
65 uniform distribution, but for all intercepts and sample sizes. The p-values for both  
66 parametric Pearson and nonparametric Pearson tests were similar and tended towards  
67 the uniform distribution with larger sample sizes.



68

69 **Figure S2.1.** Distribution of p-values for the Poisson GLMs for each dispersion test.  
 70 10,000 simulations per simulation set (intercept x sample size).

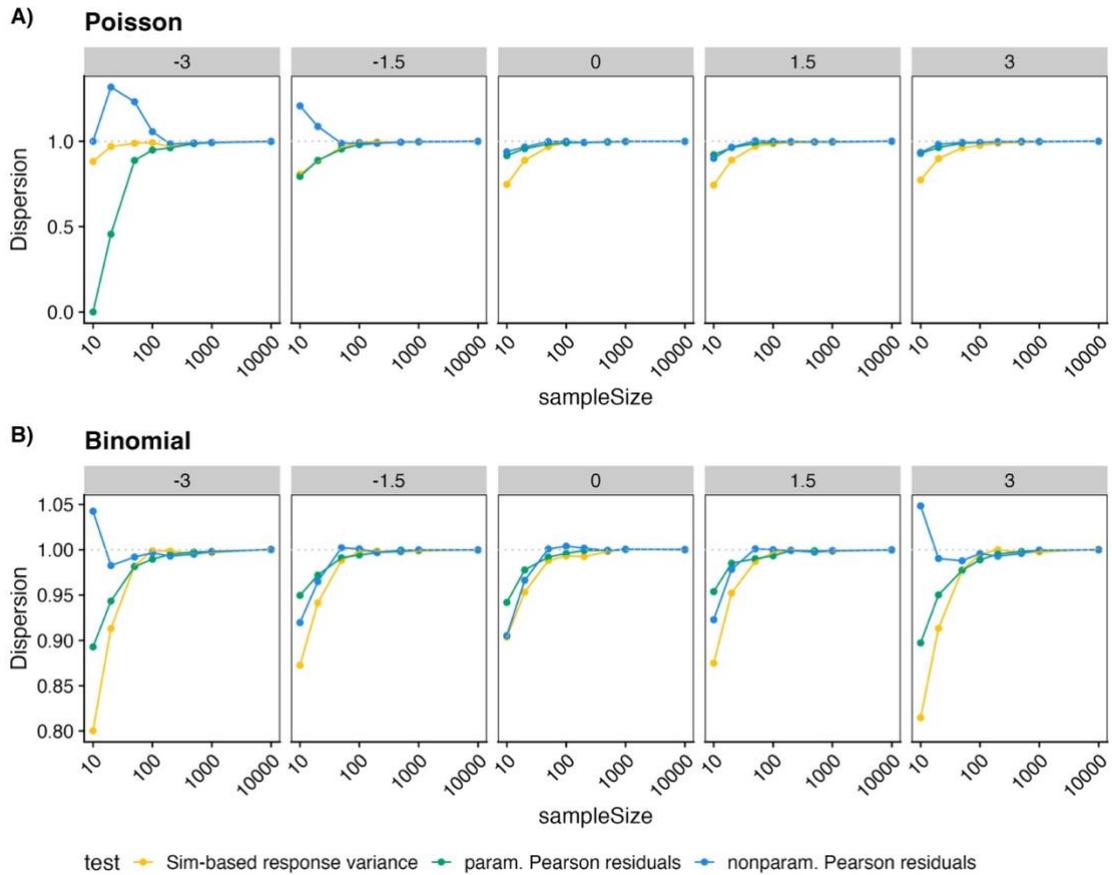


71

72 **Figure S2.2.** Distribution of p-values for the binomial GLMs for each dispersion test.  
 73 10,000 simulations per simulation set (intercept x sample size).

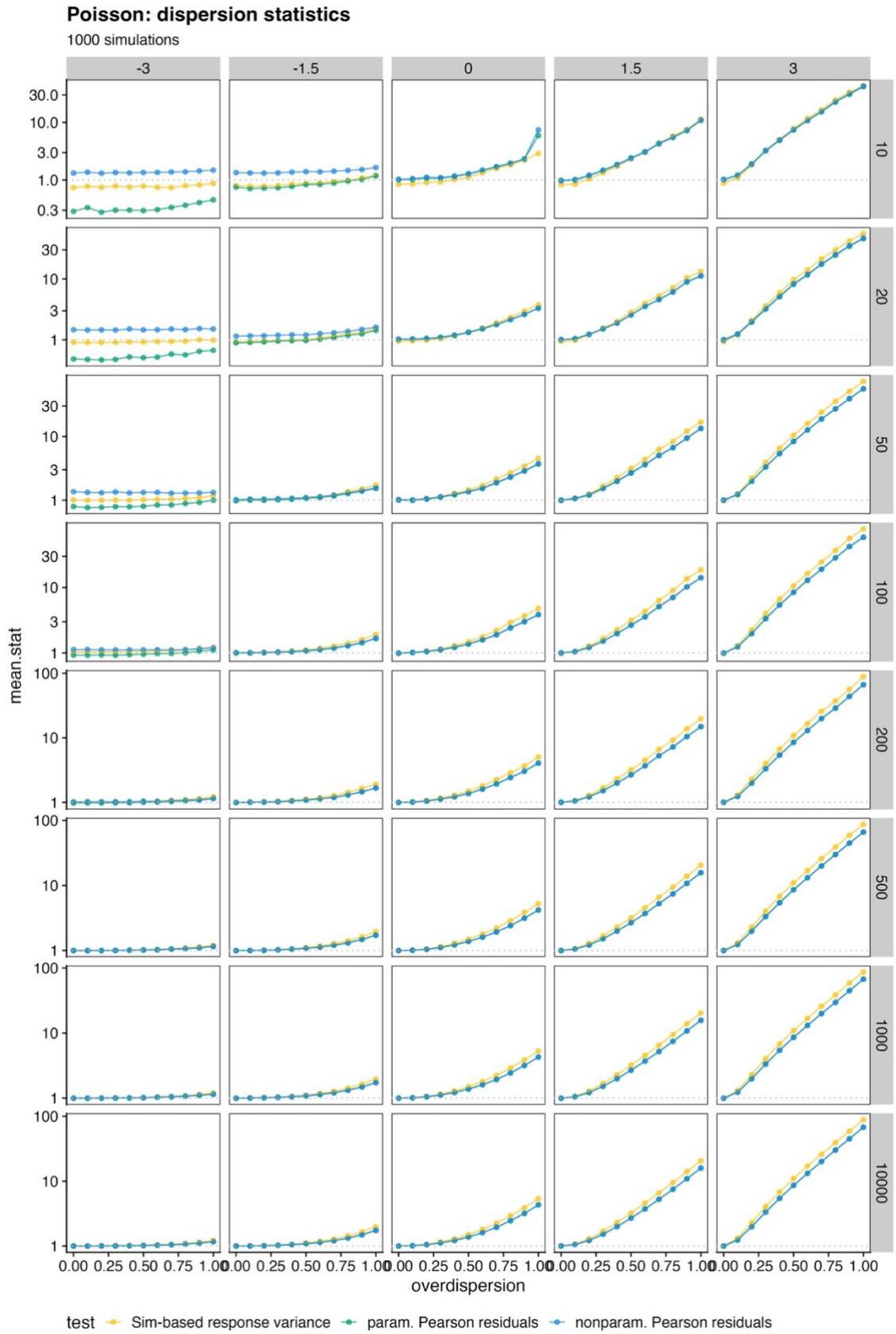
### 74 **S3. Dispersion statistics for GLMs**

75           The dispersion statistics of the tests for GLMs tended to be smaller than 1  
76 (expected value) when there was no overdispersion simulated for very small sample  
77 sizes for both binomial and Poisson distributions (Figure S3.1). The exception was the  
78 nonparametric Pearson test that presented values larger than 1 for the very small  
79 intercepts (-3 in both distributions, 3 in binomial only). When comparing dispersion  
80 statistics for the simulated overdispersed data (Figures S3.2 and S3.3), we found that  
81 both Pearson-based dispersion statistics presented similar values. In contrast, the  
82 dispersion statistic of the simulation-based response variance presented lower values for  
83 small sample sizes. The differences in dispersion statistics between tests tended to  
84 increase with the increase of simulated overdispersion, but in opposite directions for  
85 binomial and Poisson GLMs (Figure S3.4 and S3.5). Moreover, we found out that the  
86 dispersion statistics of the simulation-based response variance test depend heavily on  
87 the slope parameter of the simulated data (Figure S3.6).



88

89 **Figure S3.1.** Median of the dispersion statistics of the tests for A) Poisson and B)  
 90 binomial GLMs, simulated without overdispersion for different intercepts (panels) and  
 91 sample sizes (x-axis) for the three dispersion tests: parametric Pearson test,  
 92 nonparametric Pearson test, and simulation-based response variance test. The dotted  
 93 horizontal line indicates the ratio of 1. Values below the line are considered  
 94 underdispersion, and above the line are overdispersion. For all simulations, the slope  
 95 was fixed at 1.

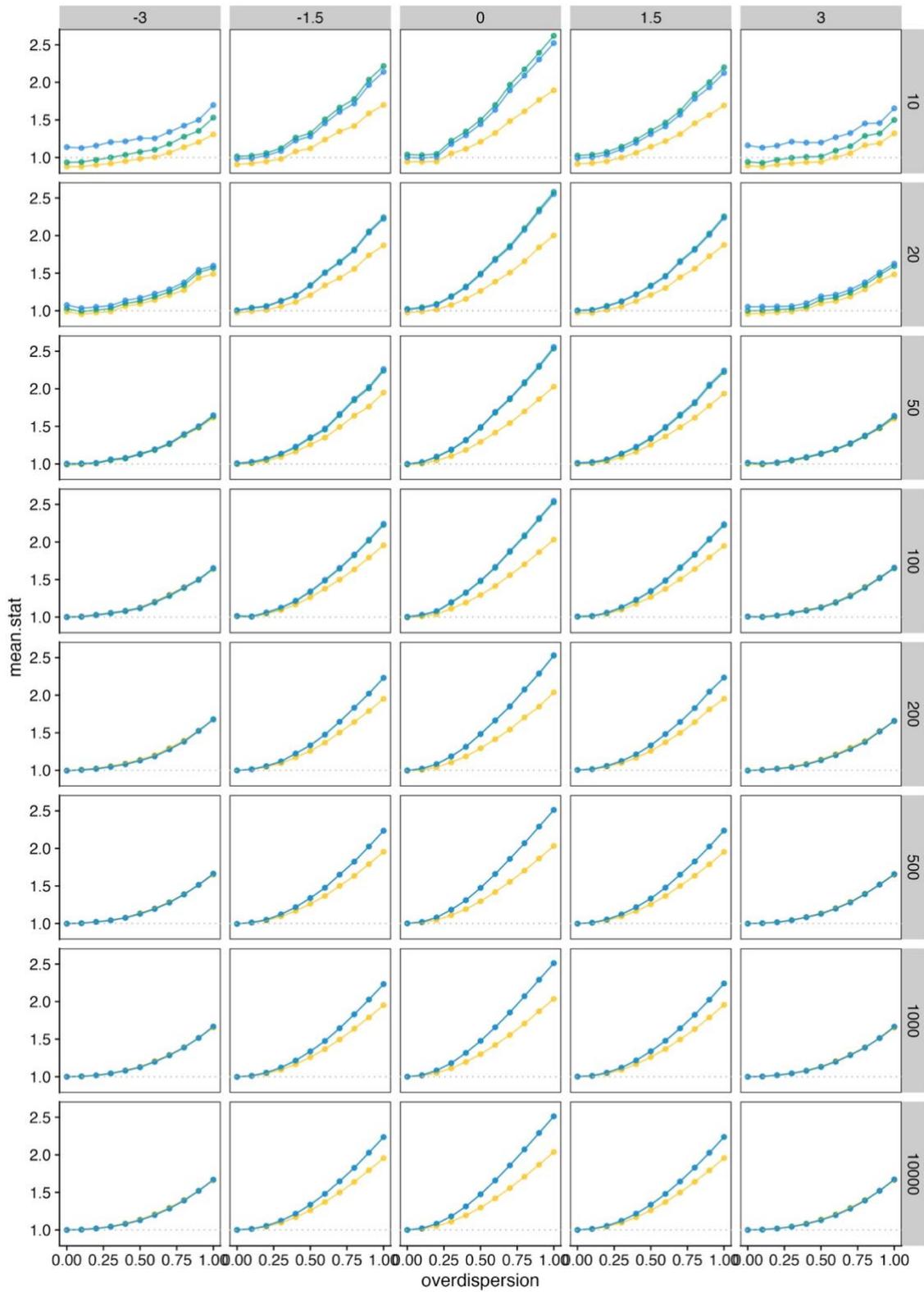


96

97 **Figure S3.2.** Dispersion statistics (median) for GLM Poisson. Notice the different y-  
98 axis scales across sample sizes.

**Binomial: dispersion statistics**

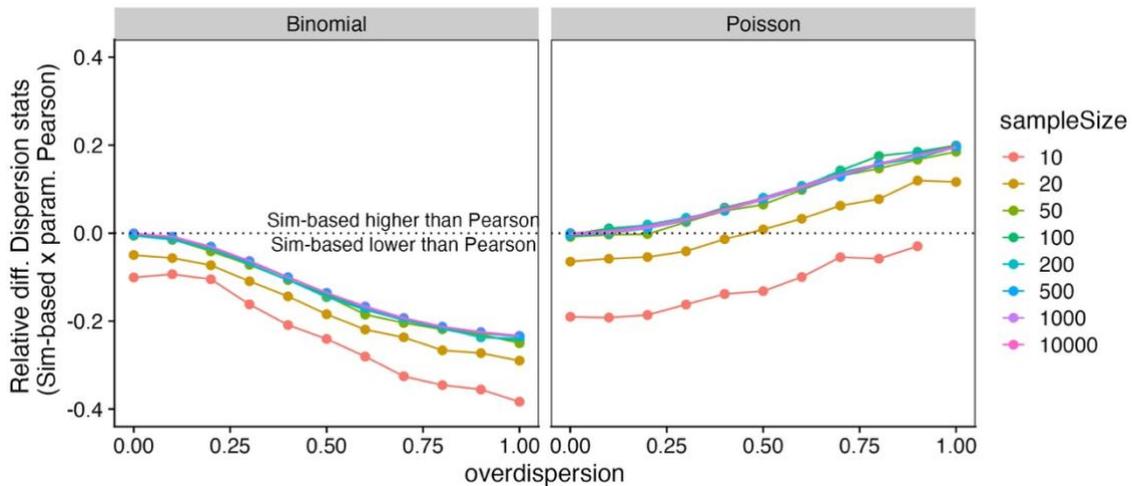
1000 sim; Ntrials=10



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

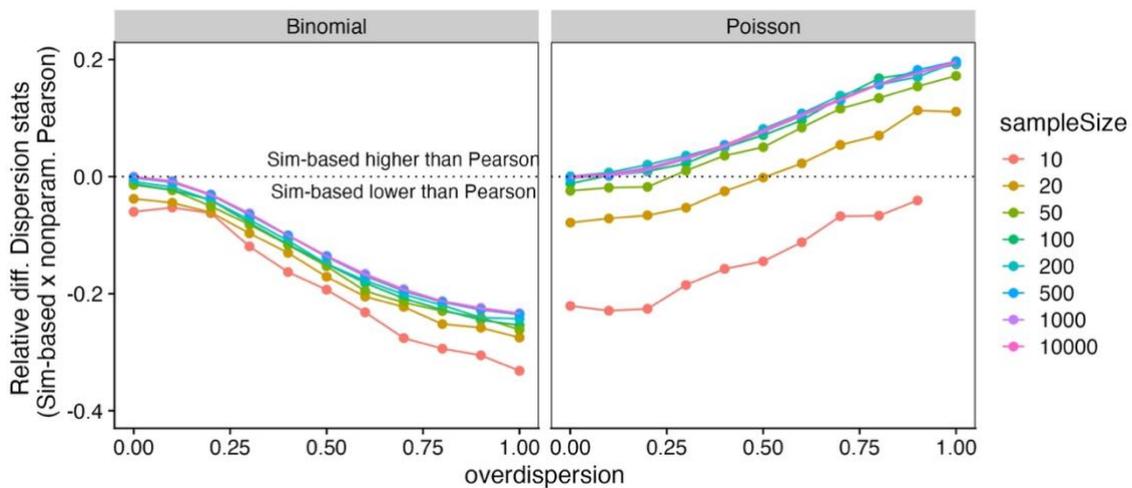
99

100 **Figure S3.3.** Dispersion statistics (median) for GLM binomial.



101

102 **Figure S3.4.** The dispersion statistics of the simulation-based response variance test are  
 103 smaller than the parametric Pearson test statistics for all binomial models and for small  
 104 sample sizes in Poisson models. The differences between the two dispersion statistics  
 105 decrease with increasing sample size (coloured lines) and increase with simulated  
 106 overdispersion in the data (x-axis). The relative differences (y-axis) were calculated by  
 107 subtracting the simulation-based dispersion statistics from the parametric Pearson  
 108 statistic, then dividing by the simulation-based statistic, and can be interpreted as the  
 109 difference in the percentage of the simulation-based statistics. The results presented are  
 110 based on 1,000 simulations with zero intercepts.

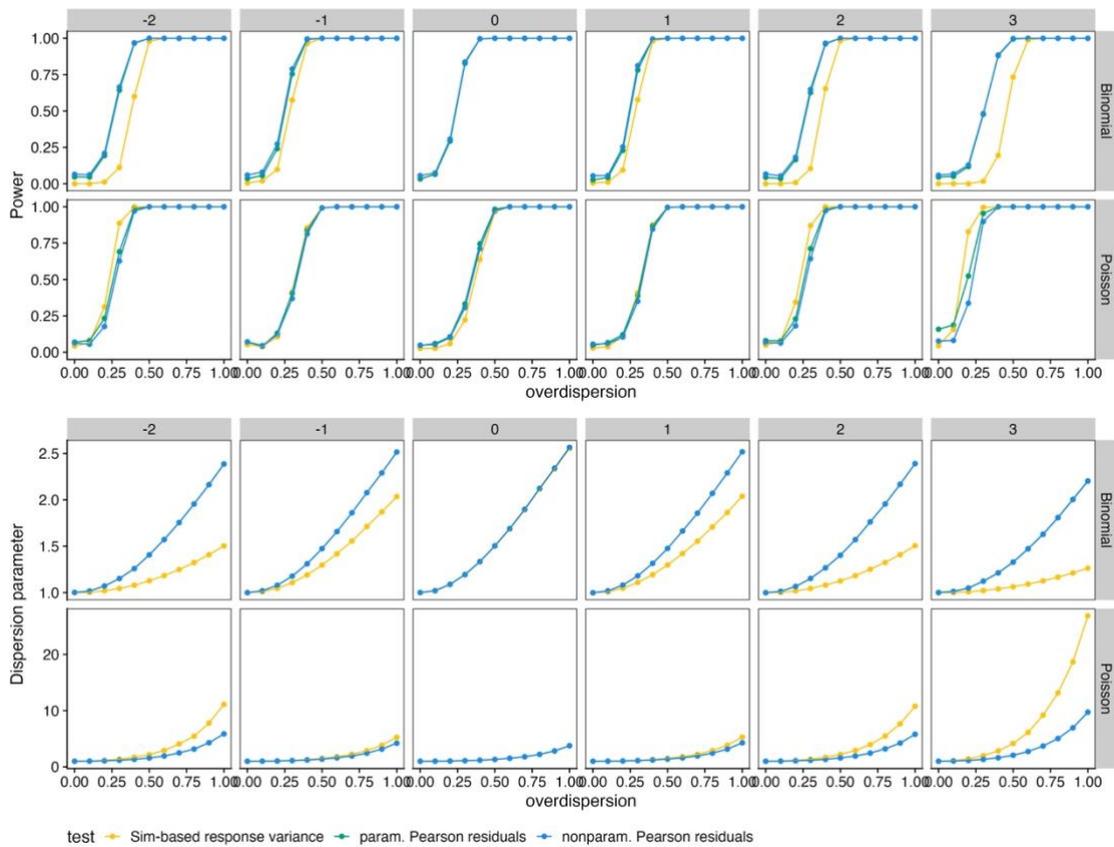


111

112 **Figure S3.5.** The dispersion statistics of the simulation-based response variance test are  
 113 smaller than nonparametric Pearson dispersion statistics for all binomial models and for  
 114 small sample sizes in Poisson models. The differences between the two dispersion  
 115 statistics decrease with increasing sample size (coloured lines) and increase with  
 116 simulated overdispersion in the data (x-axis). The relative differences (y-axis) were  
 117 calculated by subtracting the Parametric Bootstrapping statistics from the simulation-  
 118 based dispersion statistics, then dividing by the simulation-based statistics, and can be  
 119 interpreted as the difference in the percentage of the simulation-based statistics. The  
 120 results presented are based on 1,000 simulations with zero intercepts.

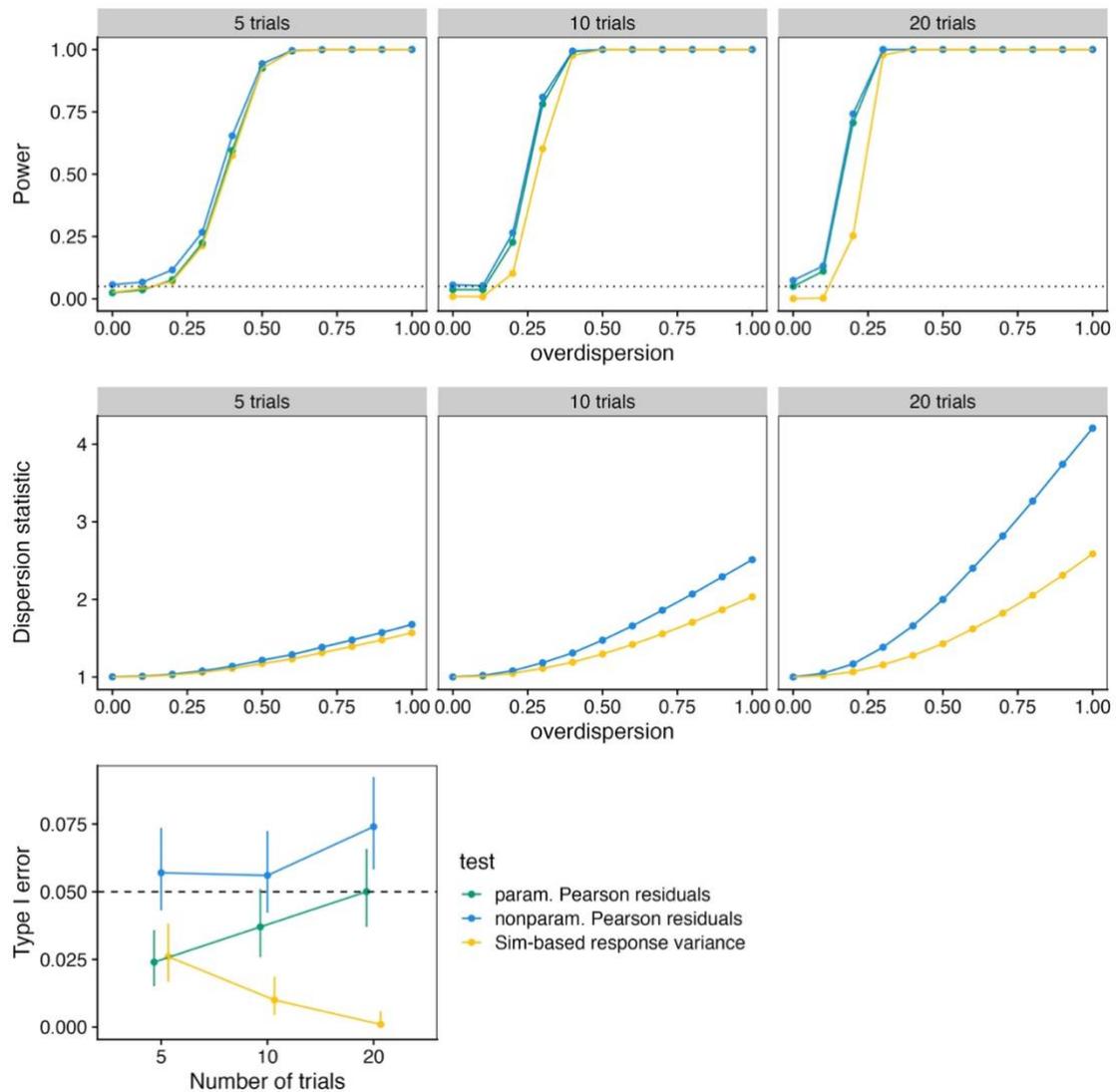
121 **S4: Expanding simulation parameters for GLMs**

122 Here, we investigated the possible influence of other parameters used to generate  
 123 the datasets for binomial and Poisson GLMs. In Figure S4.1, we investigated the power  
 124 and dispersion statistic for datasets simulated with different slopes (the default slope in  
 125 all other simulations was 1). In Figure S4.2, we investigated the effect of varying the  
 126 number of trials on the binomial GLMs in terms of power, type I error, and dispersion  
 127 statistics.



128

129 **Figure S4.1.** Power and dispersion statistics for simulations with different slopes (panel  
 130 columns) for binomial and Poisson GLMs. Number of simulations = 500; intercept = 0,  
 131 number of trials for the binomial = 10.



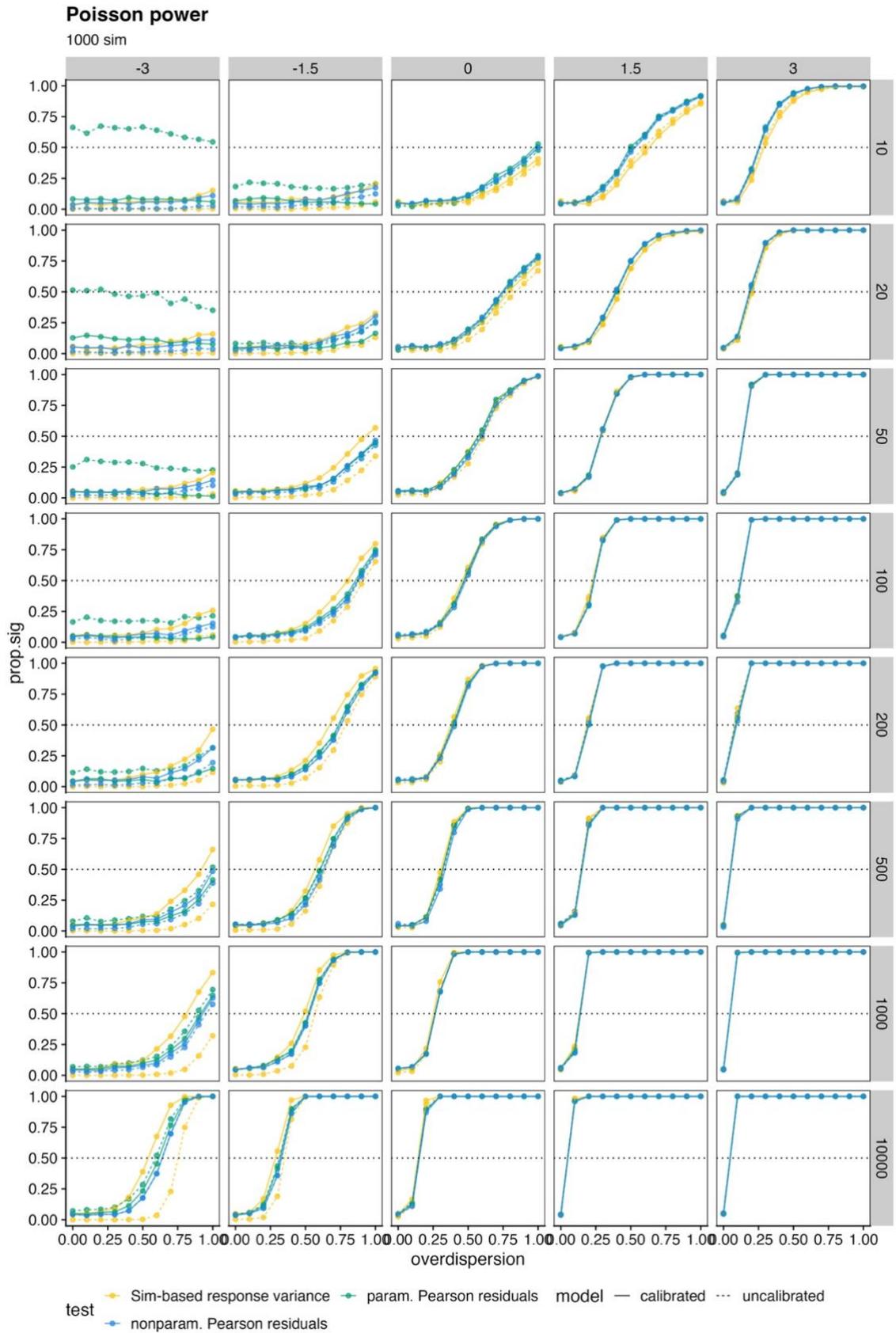
132

133 **Figure S4.2.** Power, dispersion statistics, and type I error of dispersion tests for  
 134 binomial data simulations with different numbers of trials (panel columns). The fixed  
 135 parameters are: intercept = 0, sample size = 500, slope = 1. Results for 1000  
 136 simulations.

## 137 **S5. Power for the GLMs**

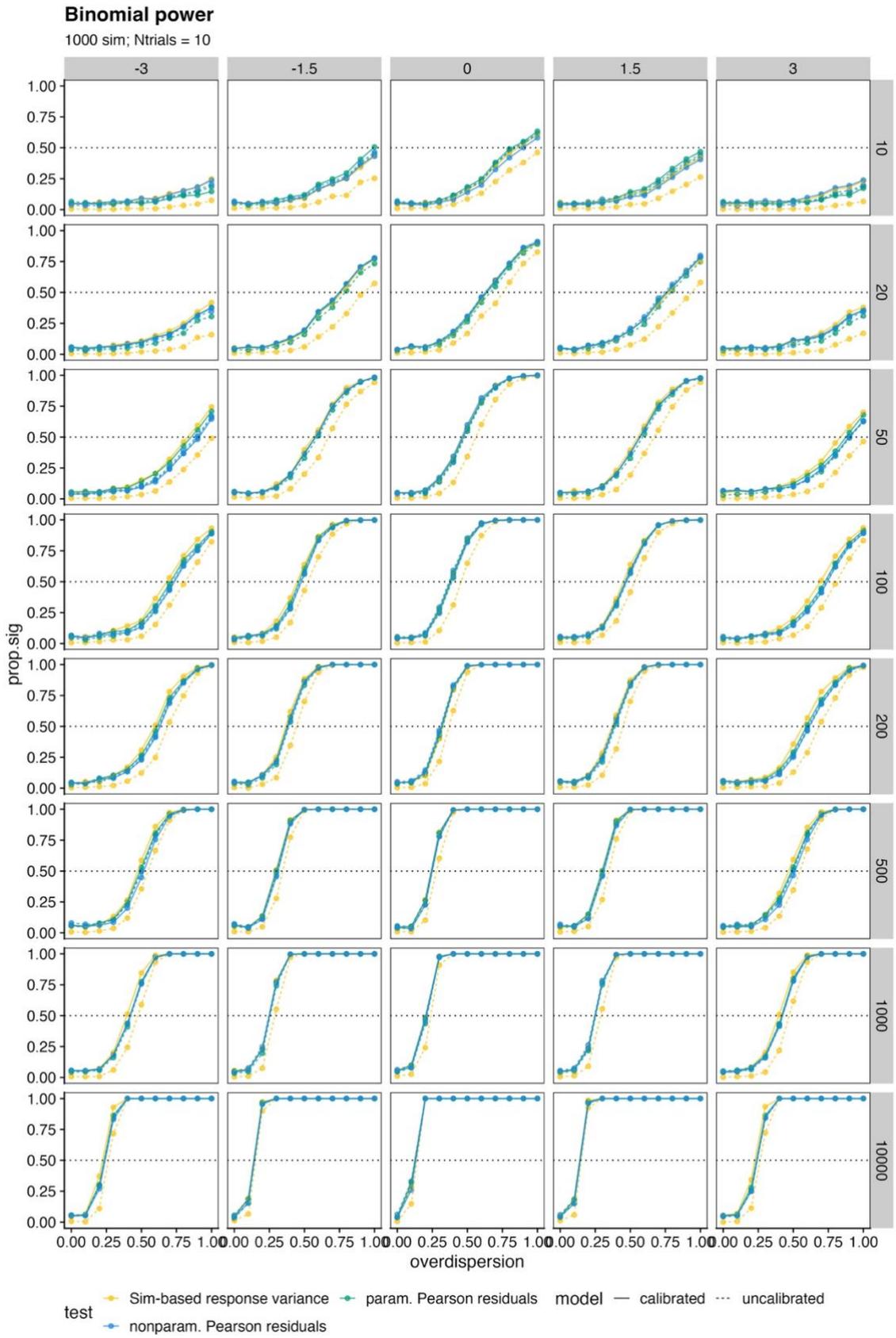
### 138 *Power calibration*

139           To investigate if the lower power of the simulation-based response variance test  
140 is a consequence of the very conservative type I error rates, we calibrated the power  
141 using the p-value at the 5% quantile of the empirical distribution of p-values where the  
142 null hypothesis was true for each set of simulations (Figures S2.1 and S2.2). This  
143 method should provide an estimate of differences in power, controlling for type I error  
144 rate (Luke et al. 2017). Figures S5.1 and S5.2 show the power (calibrated and  
145 uncalibrated) of the dispersion tests for each simulation set (intercept, sample size and  
146 overdispersion) for Poisson and binomial GLMs, respectively.



147

148 **Figure S5.1.** Power for GLM Poisson.



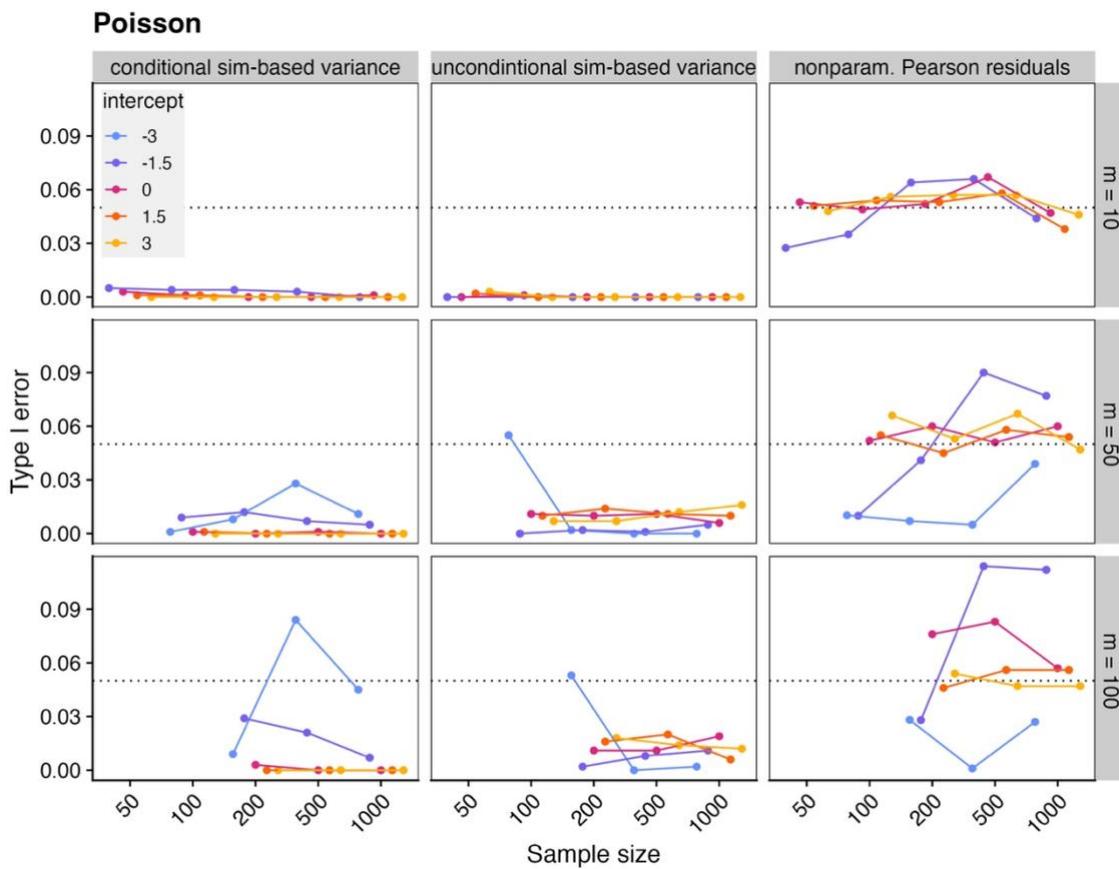
149

150 **Figure S5.2.** Power for GLM binomial.

151 **S6. Additional GLMM results**

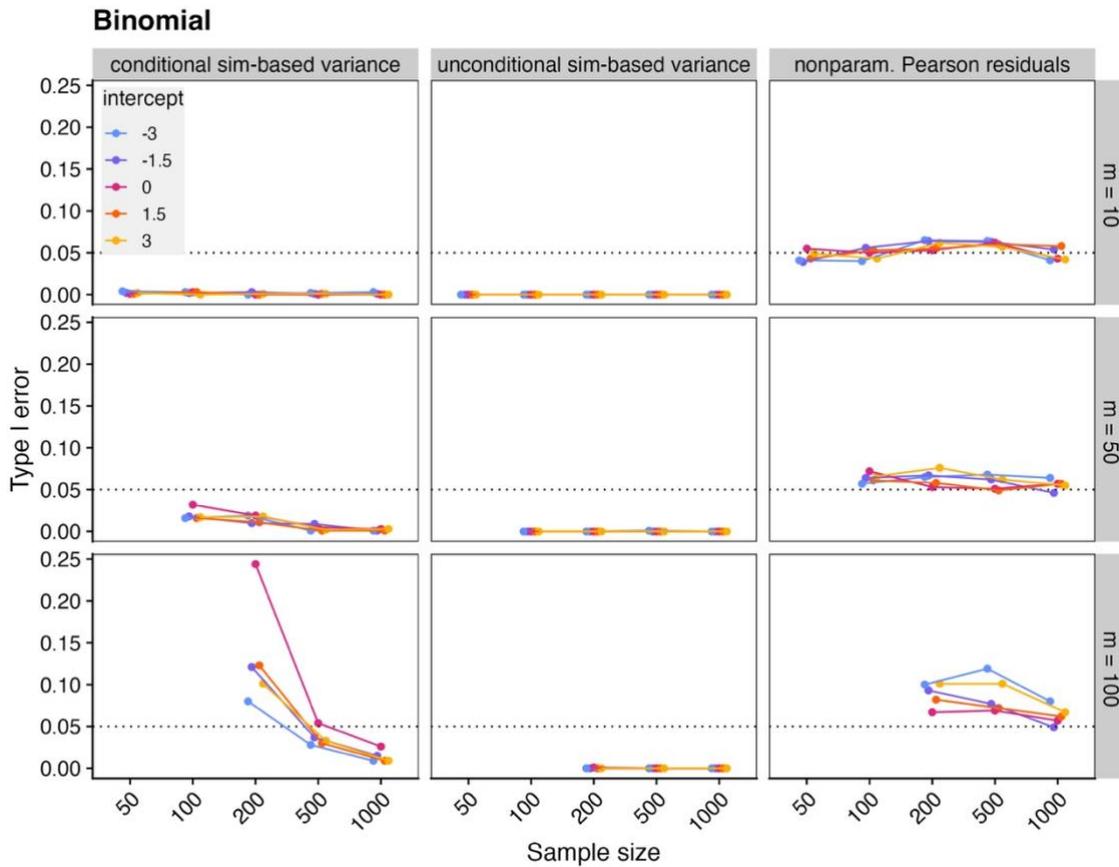
152 *Type I error rate of the alternative dispersion tests*

153 In Figures S6.1 and S6.2, we present the type I error rates for the four alternative  
 154 dispersion tests for the Poisson and binomial GLMMs, respectively, using simulated  
 155 sets of parameters: number of observations, number of groups, and intercepts.



156

157 **Figure S6.1.** Type I error rate for the three alternative dispersion tests for the Poisson  
 158 GLMMs. 1000 simulations for each parameter set. To improve visualising the different  
 159 intercept lines, the values in the x-axis were slightly displaced around the sample size  
 160 values.

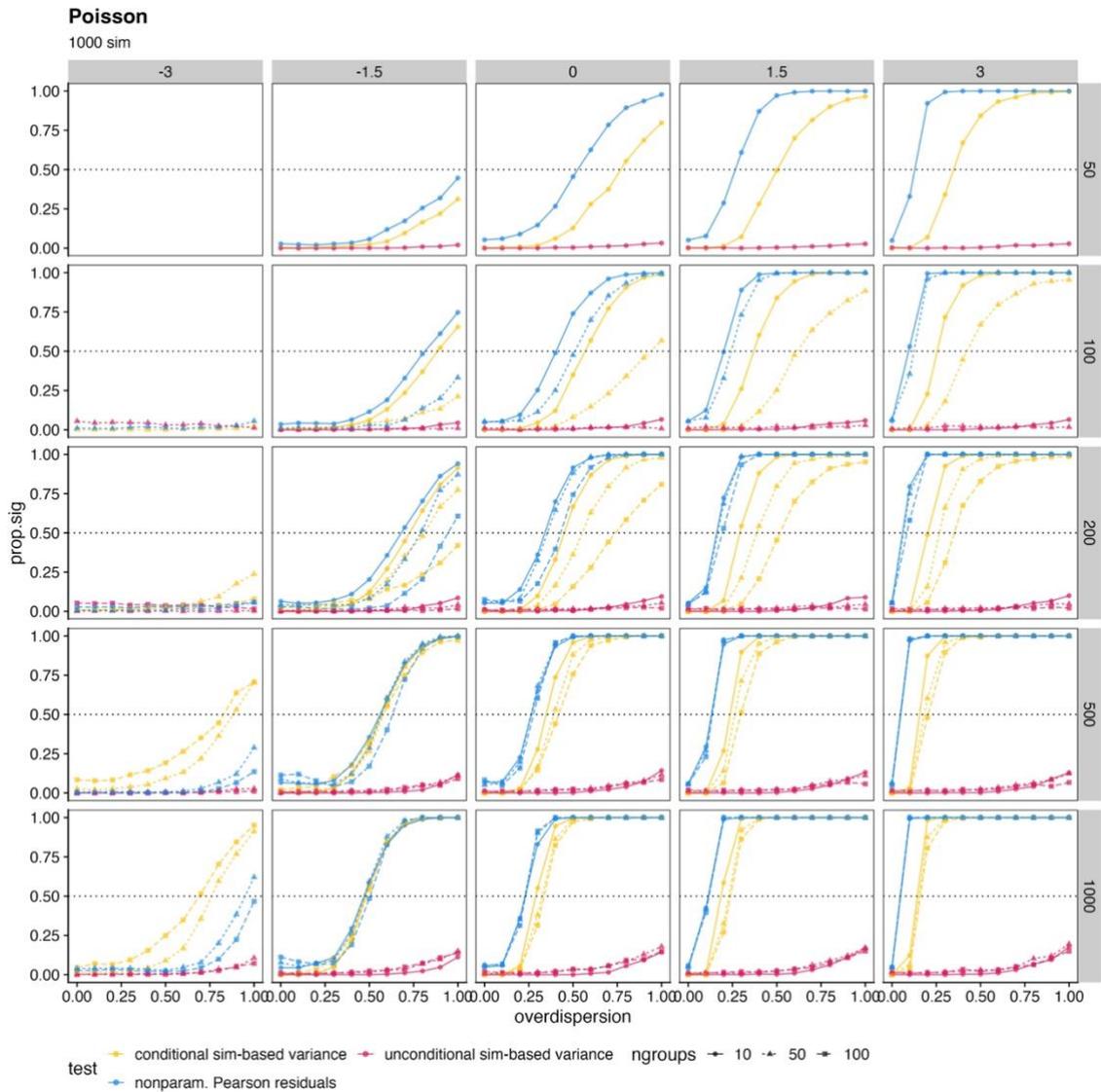


161

162 **Figure S6.2.** Type I error rate for the three alternative dispersion tests for binomial  
 163 GLMMs. 1000 simulations for each parameter set. To improve visualising the different  
 164 intercept lines, the values in the x-axis were slightly displaced around the sample size  
 165 values.

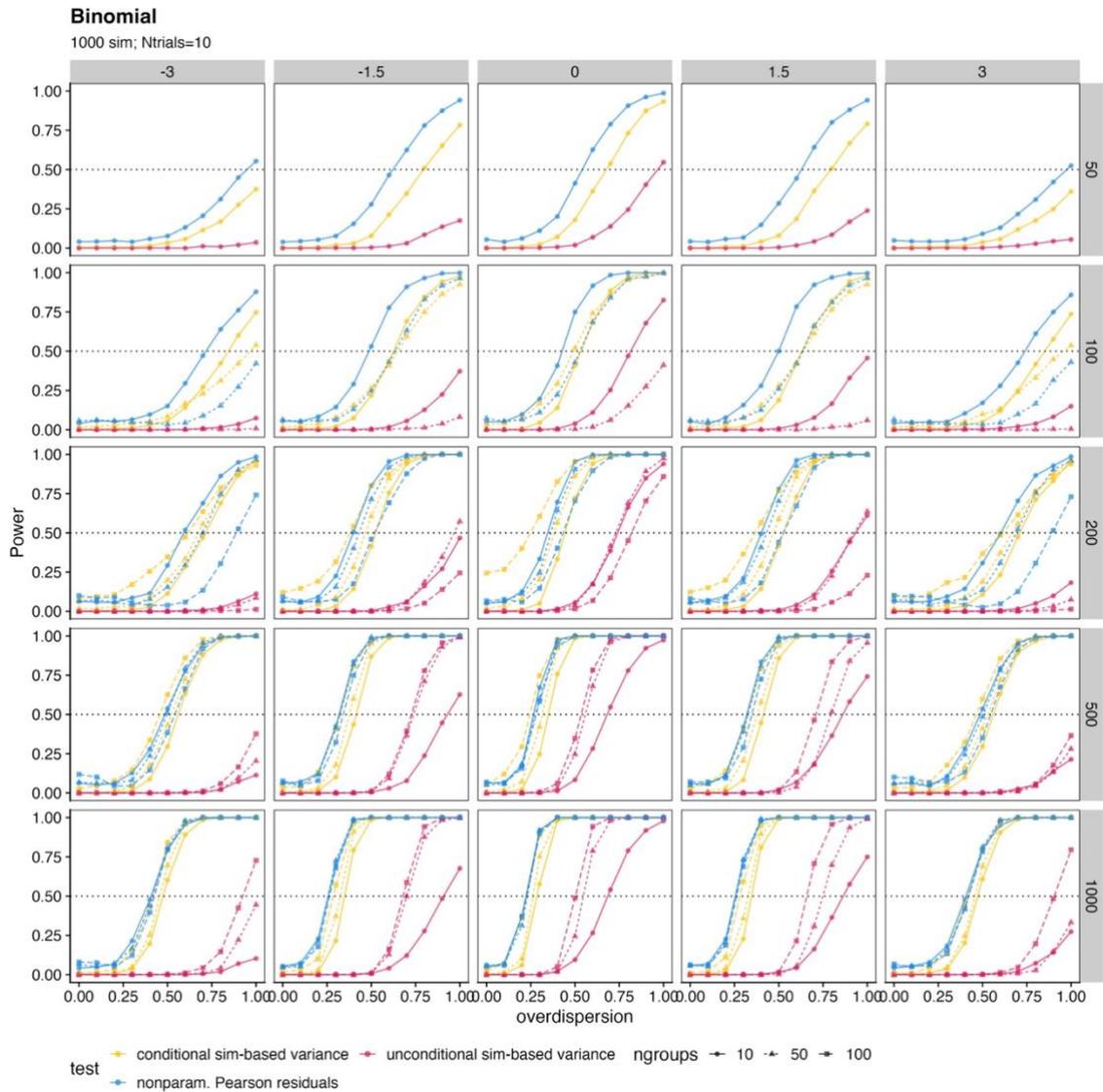
166 *Power of the alternative dispersion tests*

167 In Figures S6.3 and S6.4, we show the Power for the three alternative dispersion  
 168 tests for the Poisson and binomial GLMMs, respectively, for the simulated sets of  
 169 parameters: number of observations, number of groups, and intercepts.



170

171 **Figure S6.3.** Power of the three alternative dispersion tests for the Poisson GLMMs,  
 172 with different sample sizes (rows), intercepts (columns), and number of groups for the  
 173 random intercept (line types). The missing lines for the first panel (intercept = -3 and  
 174 sample size = 50) are due to simulation errors for some tests. For each parameter set, we  
 175 ran 1000 simulations.

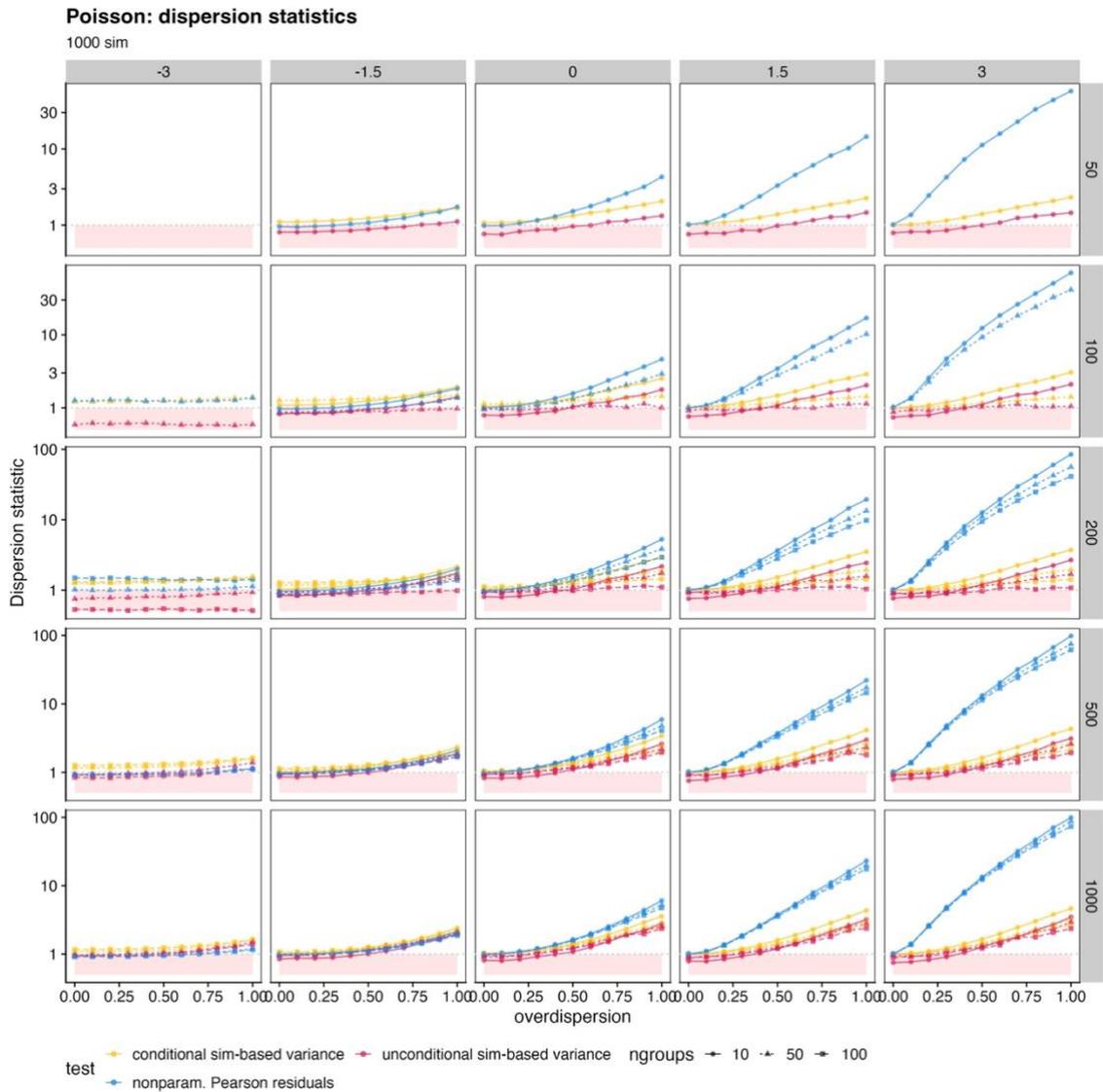


176

177 **Figure S6.4.** Power of the three alternative dispersion tests for binomial GLMMs, with  
 178 different numbers of observations (rows), intercepts (columns), and number of groups  
 179 for the random intercept (line types). 1000 simulations for each parameter set.

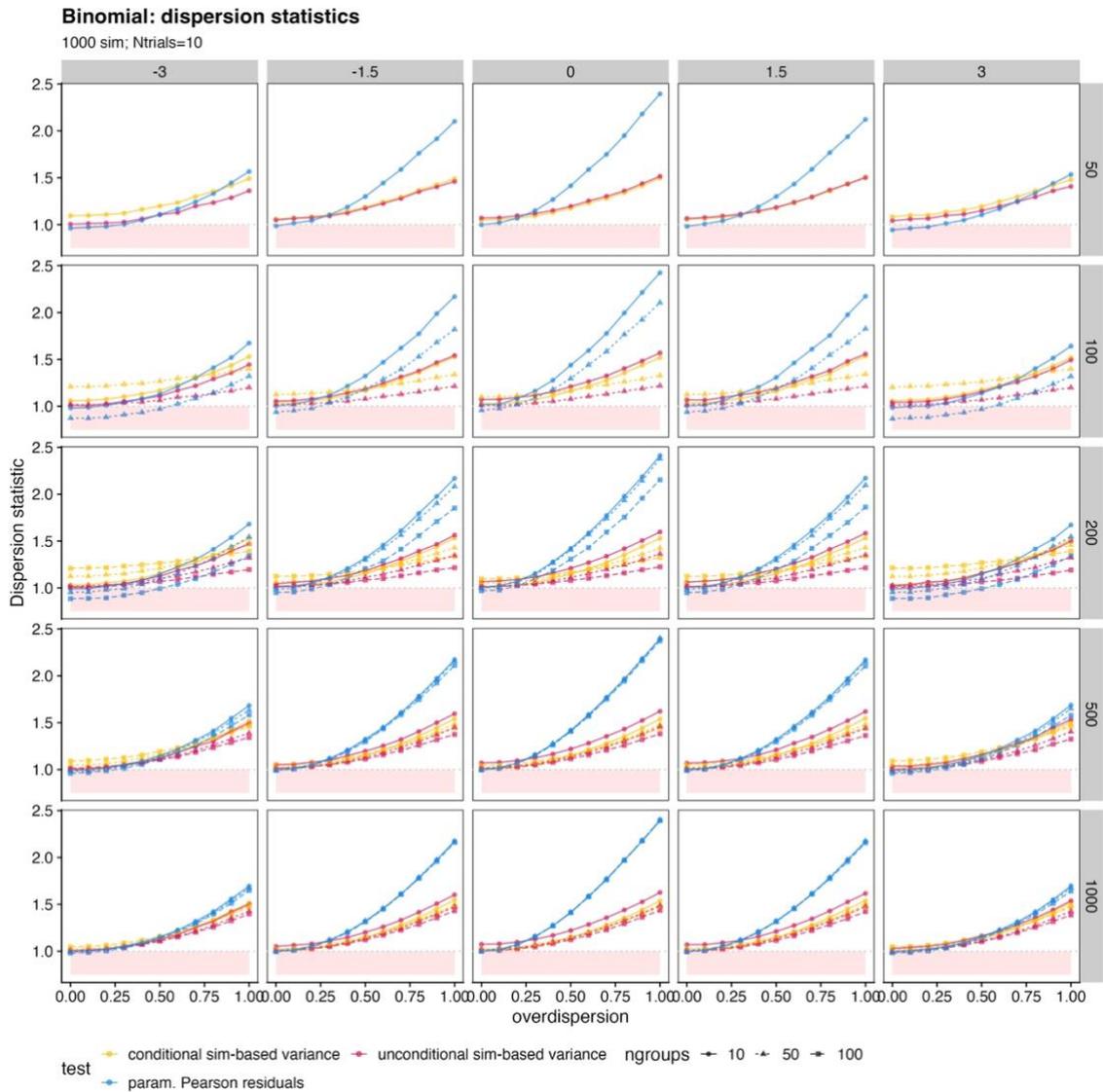
180 *Dispersion statistics of the alternative dispersion tests*

181 In Figures S6.5 and S6.6, we show the dispersion statistics for the three  
 182 alternative dispersion tests for the Poisson and binomial GLMMs, respectively, for the  
 183 simulated sets of parameters: number of observations, number of groups, and intercepts.



184

185 **Figure S6.5.** Dispersion statistics of the three alternative dispersion tests for the Poisson  
 186 GLMMs, with different numbers of observations (rows), intercepts (columns) and  
 187 number of groups for the random intercept (line types). The missing lines for the first  
 188 panel (intercept = -3 and sample size = 50 are due to simulation errors for some tests.  
 189 For each parameter set, we ran 1000 simulations.



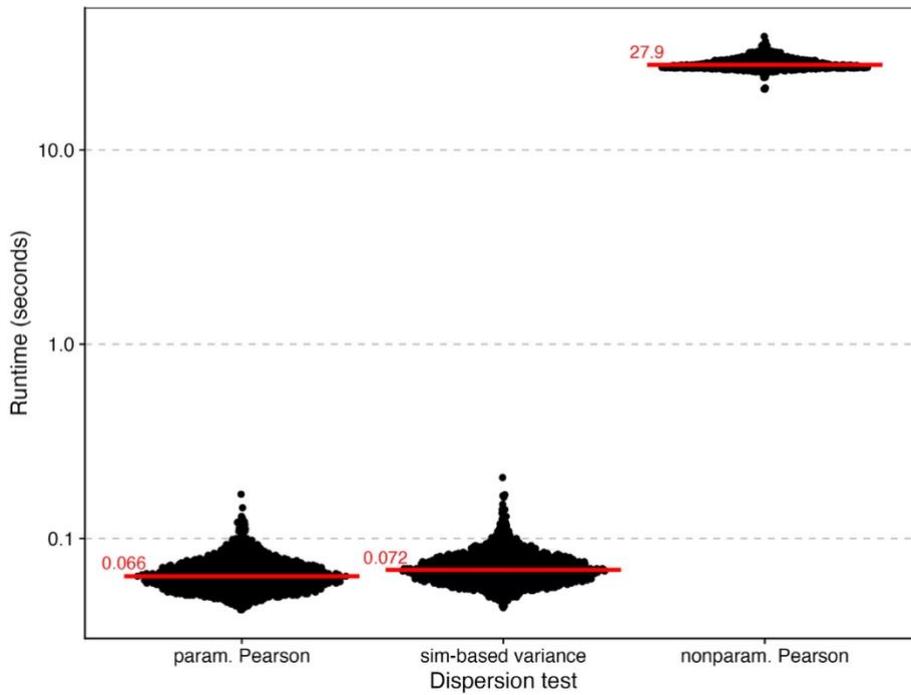
190

191 **Figure S6.6.** Dispersion statistics of the three alternative dispersion tests for binomial  
 192 GLMMs, with different numbers of observations (rows), intercepts (columns), and  
 193 number of groups for the random intercept (line types). 1000 simulations for each  
 194 parameter set.

195 *Computational runtime for tests with GLMMs*

196 We computed the run time for the three tests used for GLMMs: the parametric  
 197 Pearson test, the nonparametric Pearson test, and the simulation-based response  
 198 variance test with conditional simulations (Figure S6.7). We used 1,000 simulations of  
 199 the Poisson GLMM as an example, with an overdispersion parameter of 0.4, an  
 200 intercept of 0, a sample size of 1,000, and 100 groups. There was almost no difference

201 in computational time between the parametric Pearson test (median at 0.066 seconds)  
202 and the simulation-based response variance test (median at 0.072 seconds). As expected,  
203 the nonparametric Pearson residuals presented the largest runtime, with a median of  
204 27.9 seconds.



205

206 **Figures S6.7.** Runtime (in seconds) for each dispersion test for a Poisson GLMM  
207 simulated 1000 times with the following parameters: overdispersion parameter of 0.4,  
208 an intercept of 0, a sample size of 1,000, and a number of groups of 100. Note the y-axis  
209 at the log 10 scale.

210

## 211 **S7: Alternative simulation-based residuals dispersion test**

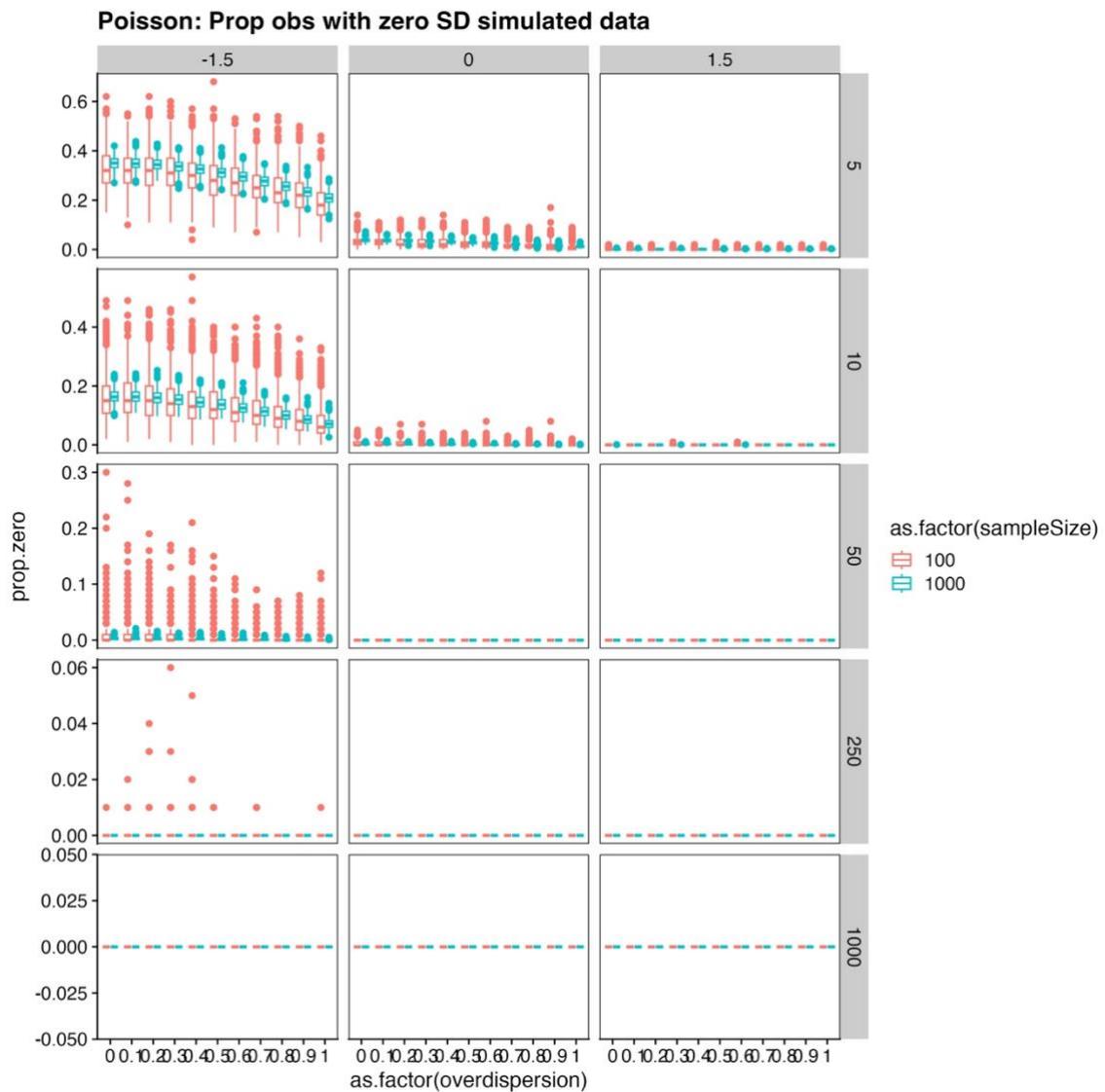
212 Another possibility for improving dispersion tests for GLMMs is to develop a  
213 simulation-based approach that shows better type I, power, and a dispersion statistic that  
214 could be interpreted similarly to the Pearson dispersion. To explore future possibilities,  
215 we briefly considered an alternative simulation-based test that attempts to approximate  
216 the Pearson residuals by dividing the observed raw residuals (observed – fitted values)  
217 by the variance of the simulated values for each observation (Equations S7.1 and S7.2).  
218 We evaluated and compared this test for Poisson and binomial GLMs and GLMMs  
219 (conditional simulations only), as we did for the other tests.

$$220 \quad \textit{Approx. Pearson observed residuals: } r_i = \frac{(y_i - \hat{\mu})}{\textit{var}(y_{is})} \quad (\textit{Equation S7.1})$$

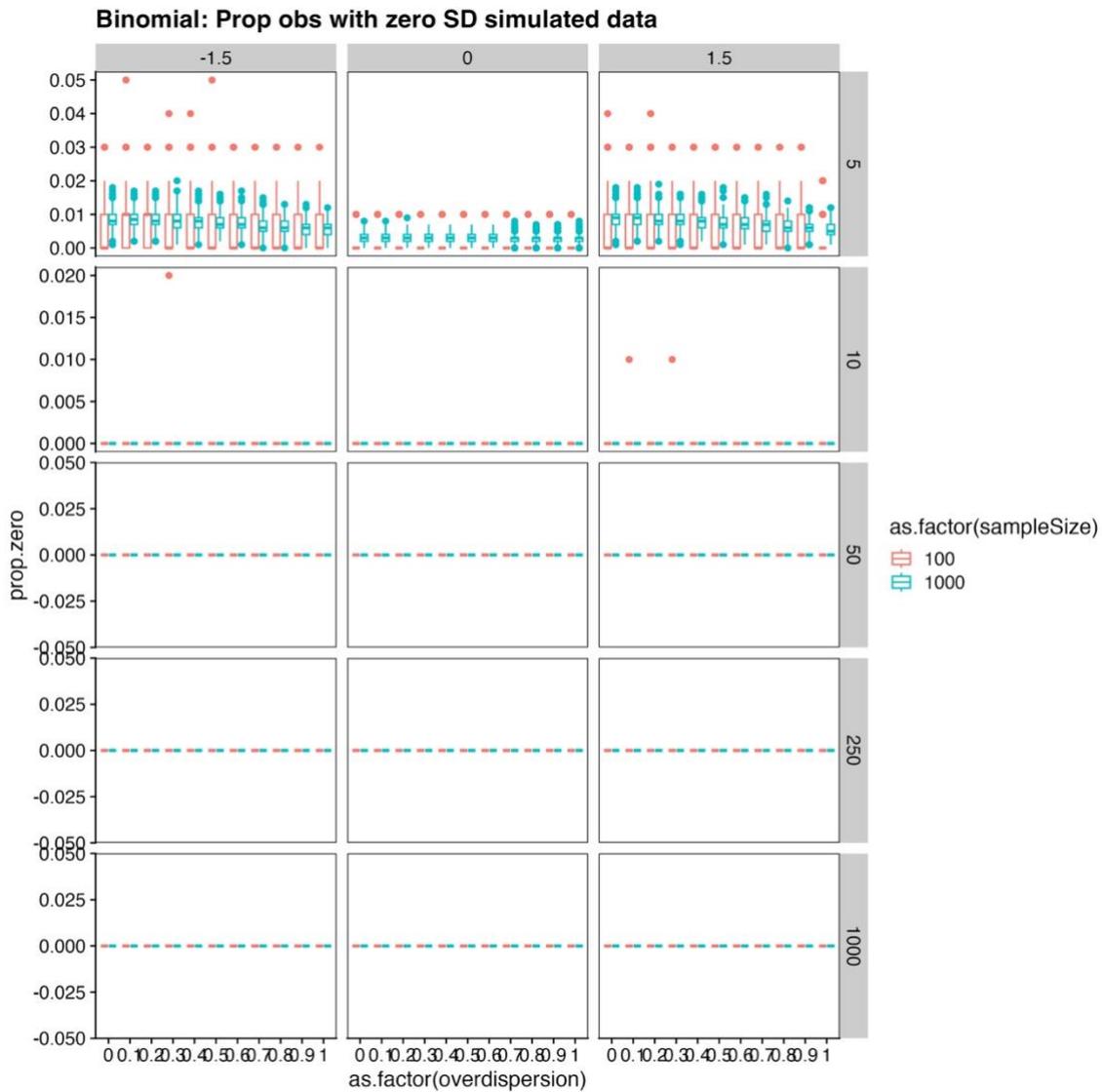
$$221 \quad \textit{Approx. Pearson simulated residuals: } r_{is} = \frac{(y_{is} - \hat{\mu})}{\textit{var}(y_{is})} \quad (\textit{Equation S7.2})$$

222 One obstacle with calculating the denominator of the approximate Pearson  
223 residuals for each observation is that the variance depends on the number of simulations  
224 and the model parameters, such as the intercept or the number of trials in the binomial  
225 GLM/GLMMs. If there are too few simulations or the intercept is very small, the  
226 chance of resulting in zero variance (all simulated values are the same) is higher for data  
227 points with small variance. To overcome this, we first evaluated the minimum number  
228 of simulations for different intercepts and sample sizes, in which all observations have  
229 estimated variances that are different from zero. For all combinations of parameters, we  
230 found that 1,000 simulations were sufficient to ensure that all variances in the simulated  
231 observations were positive (Figures S7.1 and S7.2). However, 250 simulations (the  
232 default parameter of the DHARMA package) also presented reasonable results, with the  
233 only exception being the Poisson GLMs with 30 out of 1,000 simulations (sample size

234 of 100 and intercept of -1.5) with a very low percentage of zero variances in the  
 235 simulated observations (mean of 0.01, maximum of 0.06). We are aware that the  
 236 number of zero variances in the simulations depends heavily on the simulation set, e.g.,  
 237 the number of trials for the binomial GLM. To develop an effective dispersion test, one  
 238 should consider alternatives to address this issue. For the subsequent analyses, we  
 239 excluded the simulations with zero variance in any simulated observation to compare  
 240 the alternative dispersion test with the simulation-based residuals test and the Pearson  
 241 Chi-squared dispersion test.



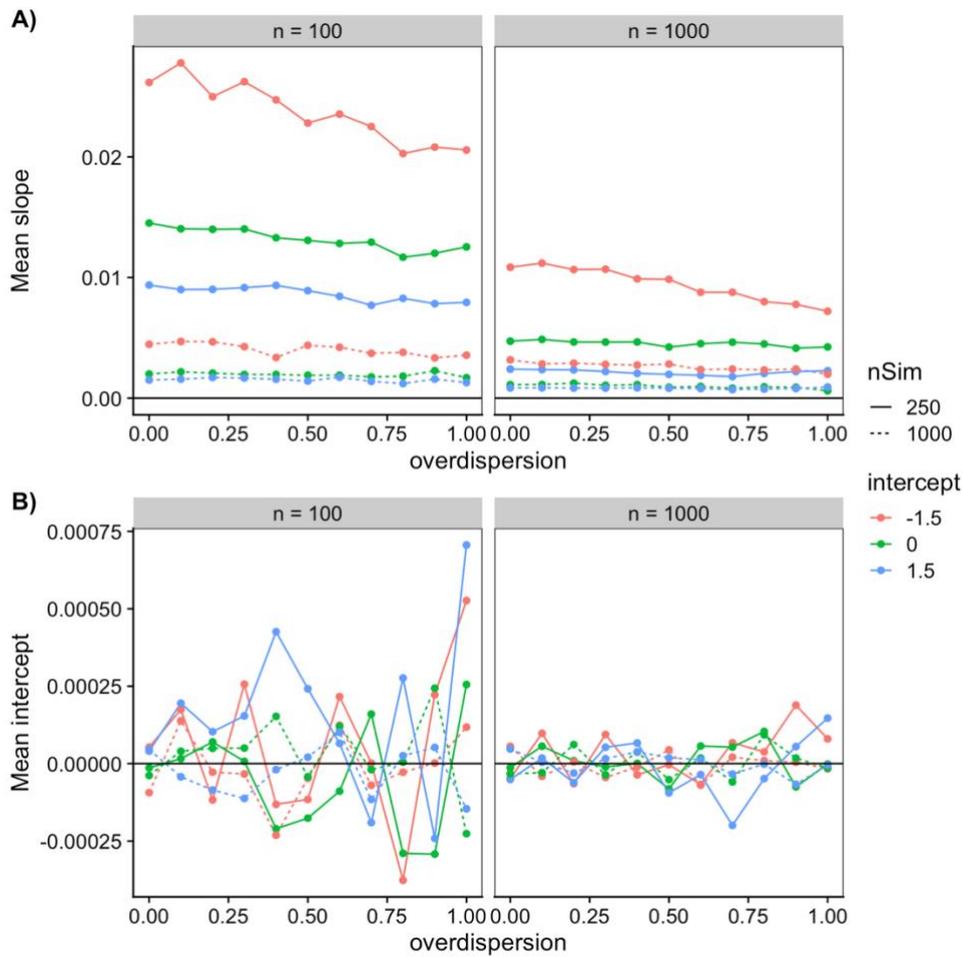
243 **Figure S7.1.** Poisson GLM: Proportion of observations with simulated zero variance in  
 244 the dataset for different combinations of intercept (columns), number of simulations  
 245 (rows) and sample sizes (colours).



246 **Figure S7.2.** Binomial GLM: Proportion of observations with simulated zero variance  
 247 in the data set for different combinations of intercept (columns), number of simulations  
 248 (rows) and sample sizes (colours). The number of trials of the binomial was set to 10 in  
 249 all simulations.  
 250

251 First, we compared the approximate Pearson residuals for GLMs with the  
 252 Pearson residuals by regressing the difference between them as the response variable  
 253 and the Pearson residuals as the predictor for the Poisson GLMs (Figure S7.3). The  
 254 intercepts for all simulation sets were nearly zero. The slope of the regression was  
 255 positive and very small for the larger number of simulations and intercepts. It means

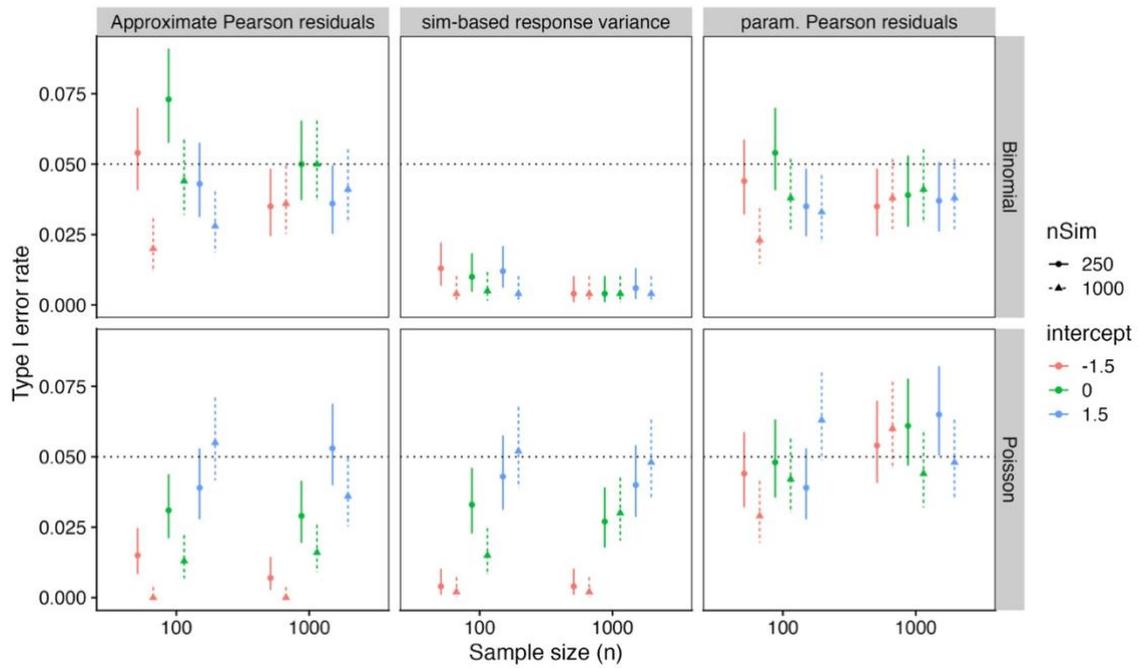
256 that the approximate Pearson tends to be slightly larger than the Pearson for larger  
 257 residuals. We did not ca



258

259 **Figure S7.3.** Mean slope (A) and intercept (B) of the regression of the difference  
 260 between the Approximate Pearson residuals and Pearson residuals as response variable  
 261 and the Pearson residuals as predictor for the Poisson GLMs.

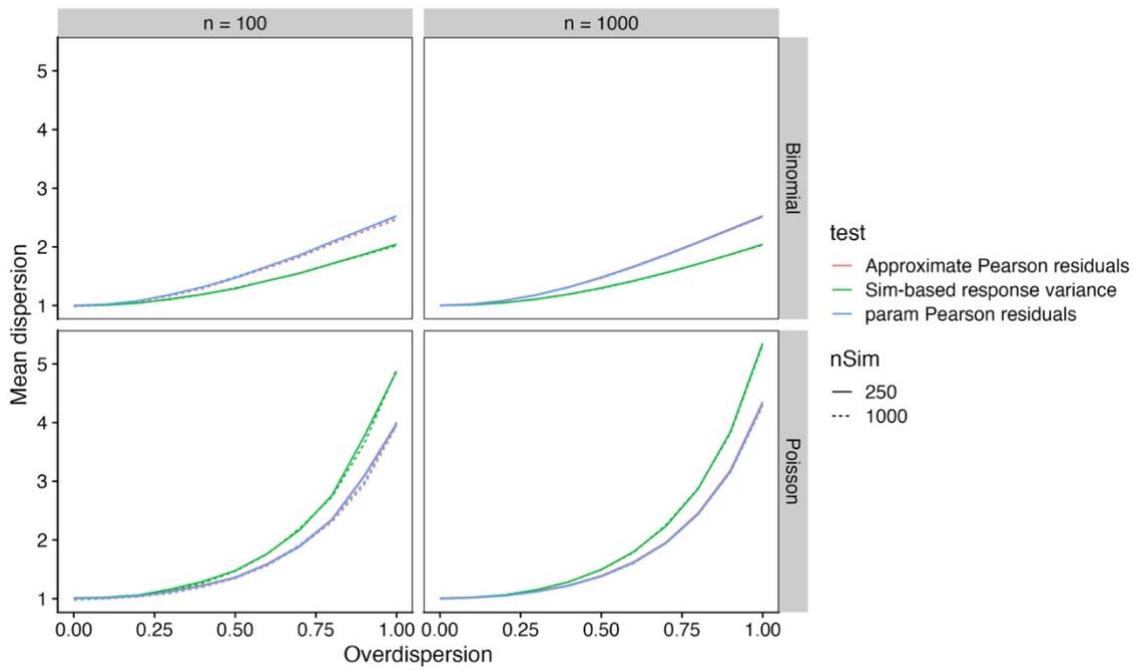
262 Type I error rates for the alternative simulation-based test, based on the  
 263 approximate Pearson residuals for GLMs, were similar to those for the simulation-based  
 264 response variance test for the Poisson model. They tended to be conservative for small  
 265 intercepts (Figure S7.4). However, for the binomial model, type I error rates were more  
 266 similar to the parametric Pearson residuals test, with values closer to 0.05 (Figure S7.4).



268

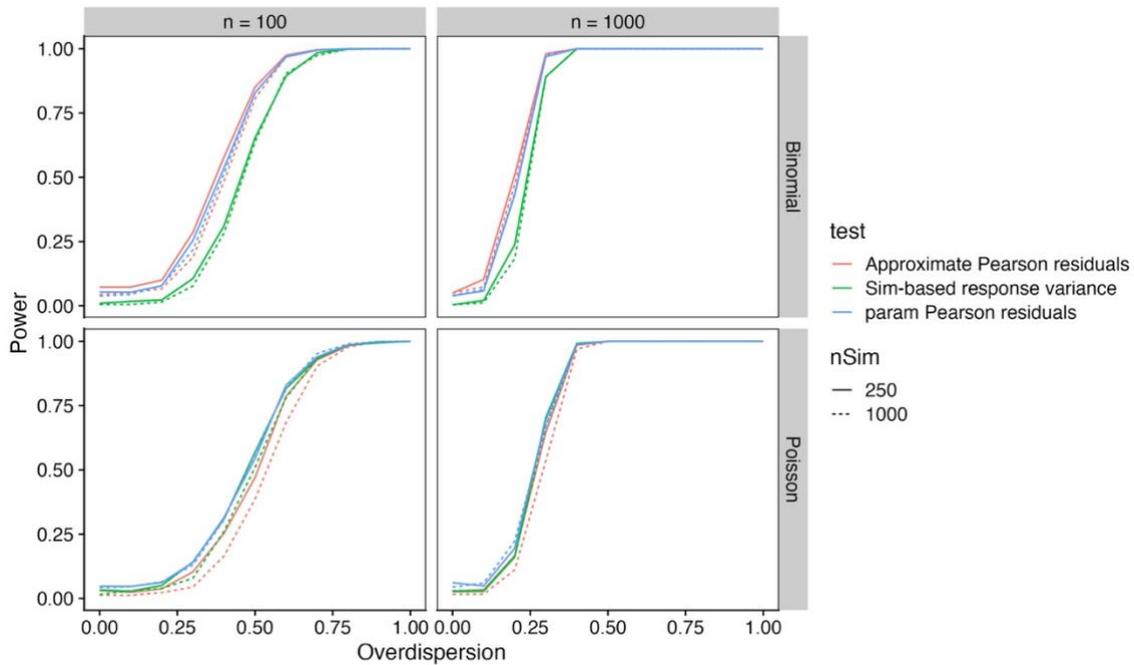
269 **Figure S7.4.** Type I error rates for GLMs comparing the parametric Pearson residuals  
 270 tests, the simulation-based response variance test and the simulation-based approximate  
 271 Pearson test.

272 The dispersion statistics for the alternative simulation-based response variance  
 273 test didn't change depending on the number of simulations and were very similar to the  
 274 parametric Pearson dispersion statistics for both GLMs (Figure S7.5). Power was very  
 275 similar among the tests for the Poisson GLM (Figure S7.6). For binomial GLMs, the  
 276 power of the alternative simulation-based residual test was high and similar to the  
 277 parametric Pearson residuals test.



278

279 **Figure S7.5.** Dispersion statistics GLMs. Simulation set with intercept = 0.

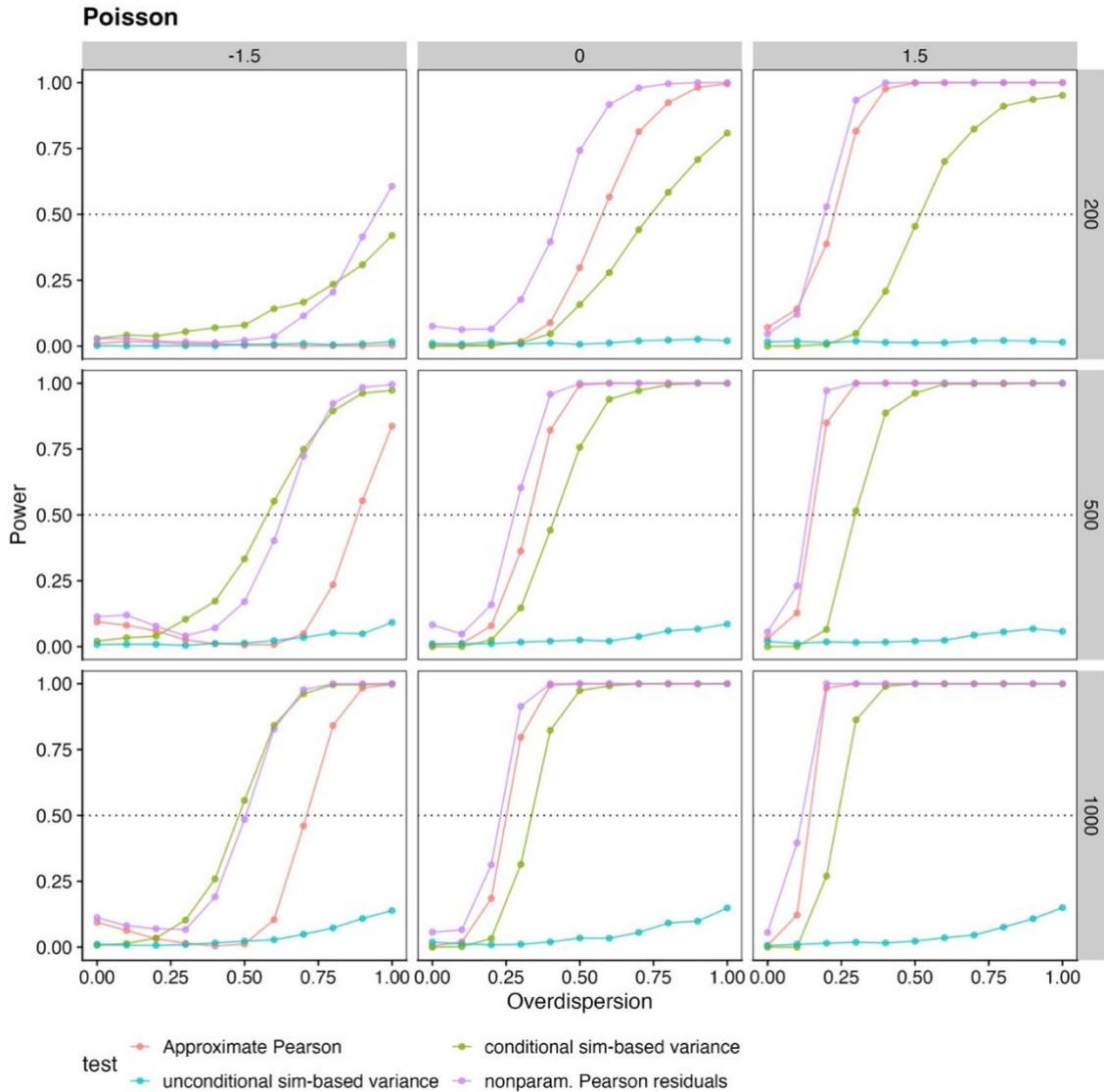


280

281 **Figure S7.6.** Power GLMs. Simulation set with intercept = 0.

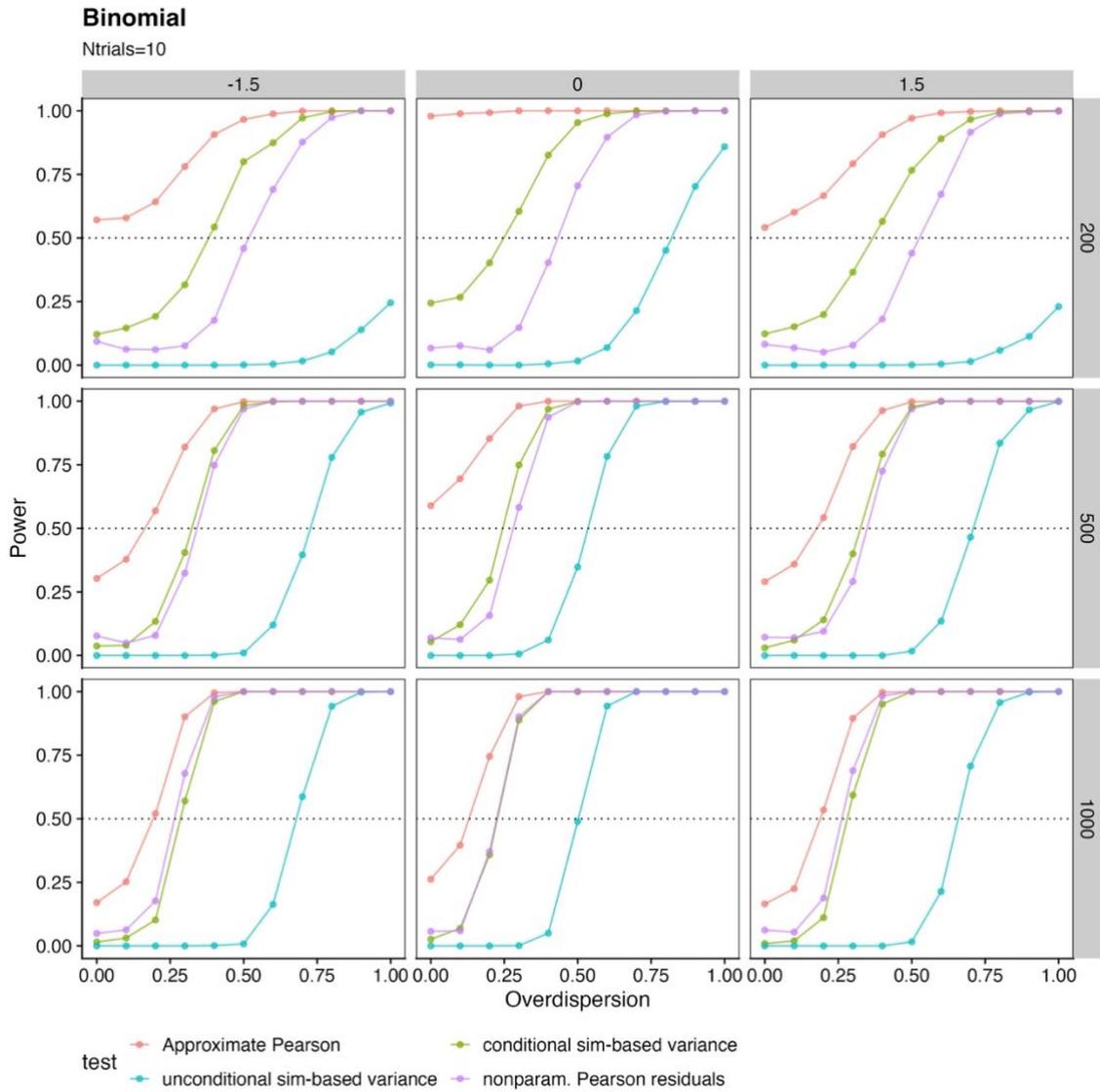
282 For the GLMM simulations, we fixed the number of groups at 100 and the  
 283 number of simulations at 250 to compare with the cases where the Pearson Chi-squared  
 284 test fails. We compared sample sizes of 200, 500, and 1000 observations and intercepts

285 of -1.5, 0, and 1.5. We excluded simulations with zero variance in the simulated  
 286 observations (specifically, for Poisson GLMMs, which accounted for less than 0.1% of  
 287 the simulations). For GLMMs, we used only the conditional simulations, which have  
 288 been proven to yield better dispersion test results.



289

290 **Fig S7.7.** Power for Poisson GLMMs for the alternative simulation-based test using an  
 291 approximation for Pearson residuals compared with the other tests assessed in the study.  
 292 1000 simulations for each parameter set: intercept (panel columns) and sample size  
 293 (panel rows). The fixed parameters are slope = 1, number of groups = 100, and random  
 294 effects variance = 1.



295

296 **Fig S7.8.** Power for binomial GLMMs for the alternative simulation-based test using an  
 297 approximation for Pearson residuals compared with the other tests assessed in the study.  
 298 1000 simulations for each parameter set: intercept (panel columns) and sample size  
 299 (panel rows). The fixed parameters are slope = 1, number of groups = 100, random  
 300 effects variance = 1, number of trials = 10.

301

302

303 **S8. Parametric Pearson test with approximated residual degrees of**  
304 **freedom for GLMMs**

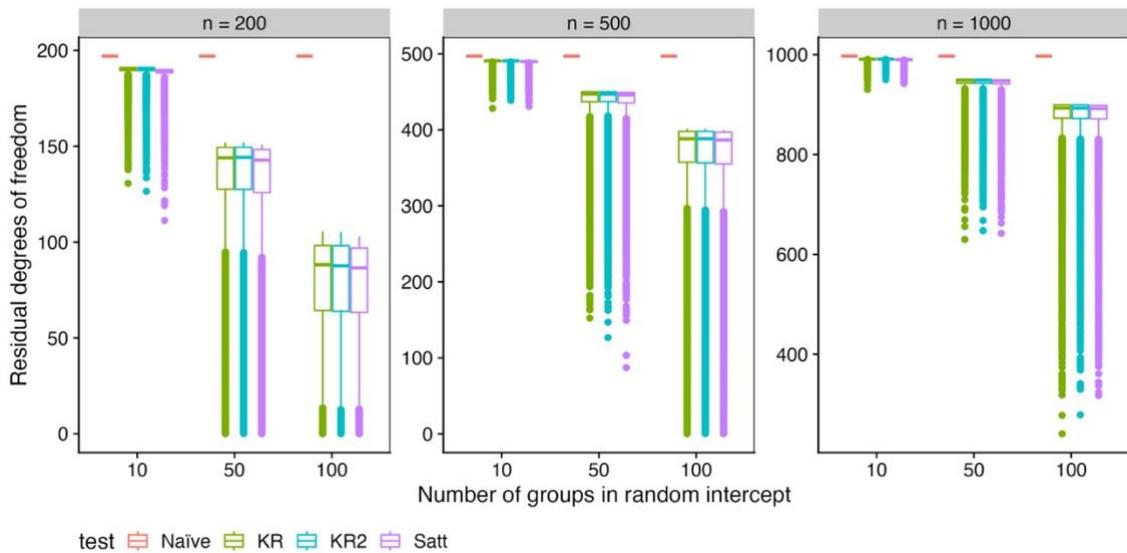
305 Degrees of freedom (*df*) are not always known for GLMMs with complex  
306 hierarchical structures and limit the use of the parametric Pearson test because it  
307 depends on it for evaluating overdispersion with the Chi-squared distribution.  
308 Moreover, our results show that using the naïve *df* is problematic for testing dispersion  
309 when you have a large number of groups in the random intercept. The two most  
310 suggested methods to approximate *df* of mixed-effect models, the Satterthwaite (1946)  
311 and the Kenward-Roger (Kenward & Roger 2009), were developed for LMMs to  
312 account for the effect of the covariance structure on *df* and standard errors. Stroup et al.  
313 (2013) suggested that the adjustment is also accurate for GLMMs. However, none of the  
314 most used R packages use any correction for the degrees of freedom for GLMMs. The  
315 few R packages that provide those approximations, e.g. *lmerTest* (Kuznetsova et al.,  
316 2017; Kuznetsova et al., 2020) that relies on *pbkrtest* (Halekoh & Højsgaard 2014), are  
317 only implemented for LMMs.

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

- 322 - The Kenward-Roger (KR) bias-corrected variance-covariance matrix for the  
323 fixed effect parameters and degrees of freedom from Kenward & Roger (1997).
- 324 - The improved correction of the Kenward-Roger (KR2) returns an improved  
325 correction given in Kenward & Roger (2009).
- 326 - The Satterthwaite correction (Sat) from Satterthwaite (1946).

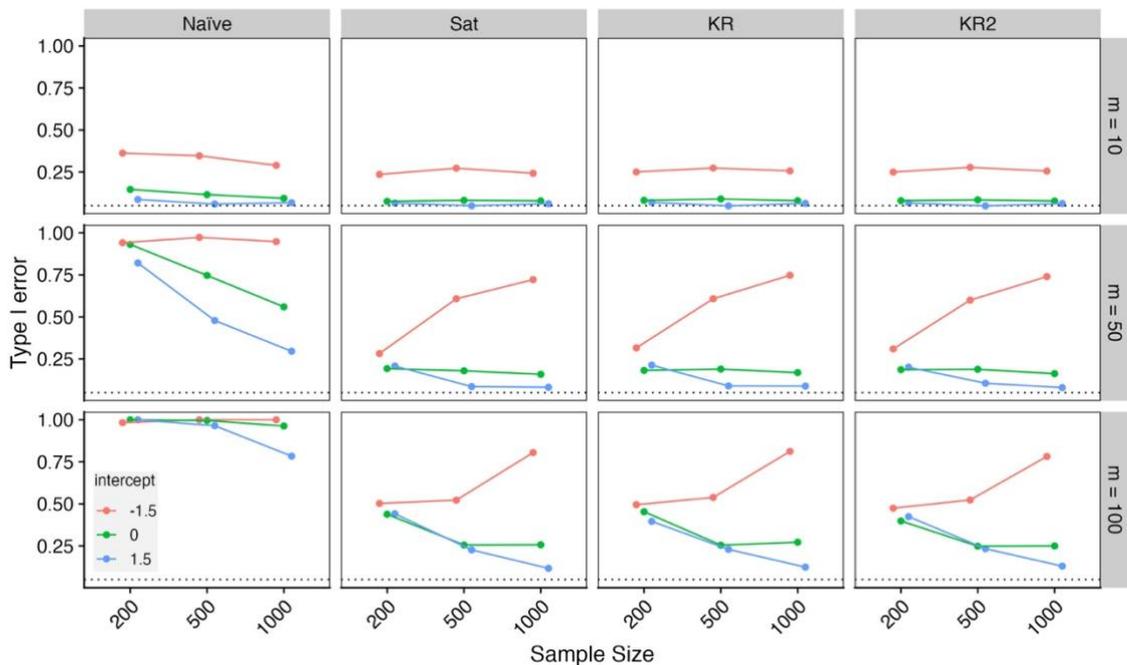
327           Our test results show that all three correction methods presented very similar  
328 residual  $df$  for all simulation settings (Figure S8.1), which resulted also in very similar  
329 test results (e.g., Figure S8.2 for type I error). Given the high similarity among tests for  
330 the different residual  $df$  corrections, we show and discuss the results for the KR2 test in  
331 comparison with the parametric Pearson “naïve” test and the alternative GLMM tests  
332 (nonparametric Pearson and simulation-based response variance test with conditional  
333 simulations). In Figure S8.3, we observe that the correction for the residual  $df$  corrected  
334 the dispersion statistics towards 1 for simulations without overdispersion, except for the  
335 very small intercept (-1.5). This results in the two-sided dispersion test being less prone  
336 to being significant, given the very low dispersion parameter (detecting underdispersion  
337 instead of overdispersion).

338           Although the parametric Pearson tests with the approximated residual degrees of  
339 freedom performed much better than those with the “naïve” residual  $df$ , they still  
340 underperformed compared to the nonparametric version when having a large number of  
341 groups in the random effects (Figure S8.4), especially for very small intercepts and  
342 sample sizes.



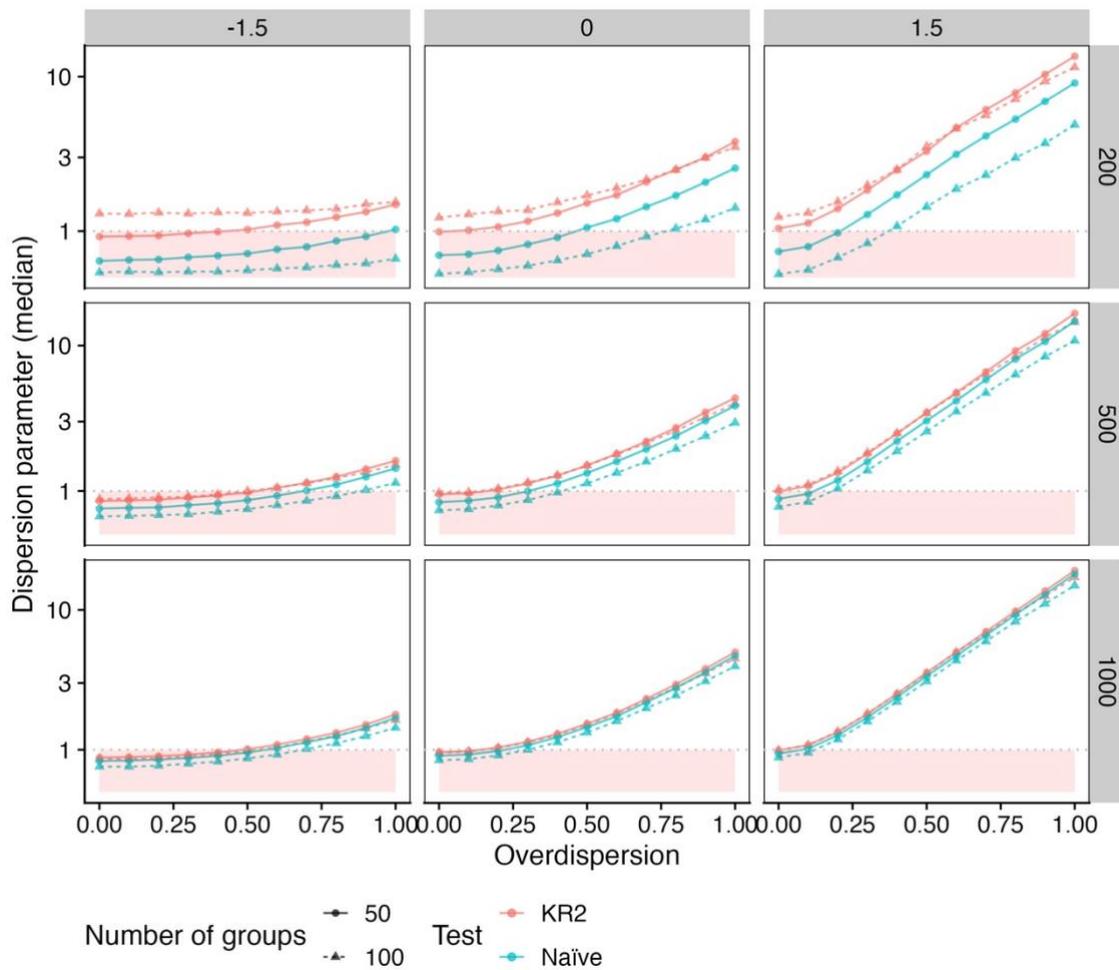
343

344 **Figure S8.1.** Residual degrees of freedom for the different correction methods for  
 345 Poisson GLMMs with different numbers of groups in the random intercept (x-axis) and  
 346 sample sizes (panel columns). Please refer to the main text above to relate to each  
 347 applied correction. 1,000 simulations for each parameter setting, slope = 1, random  
 348 intercept variance = 1.



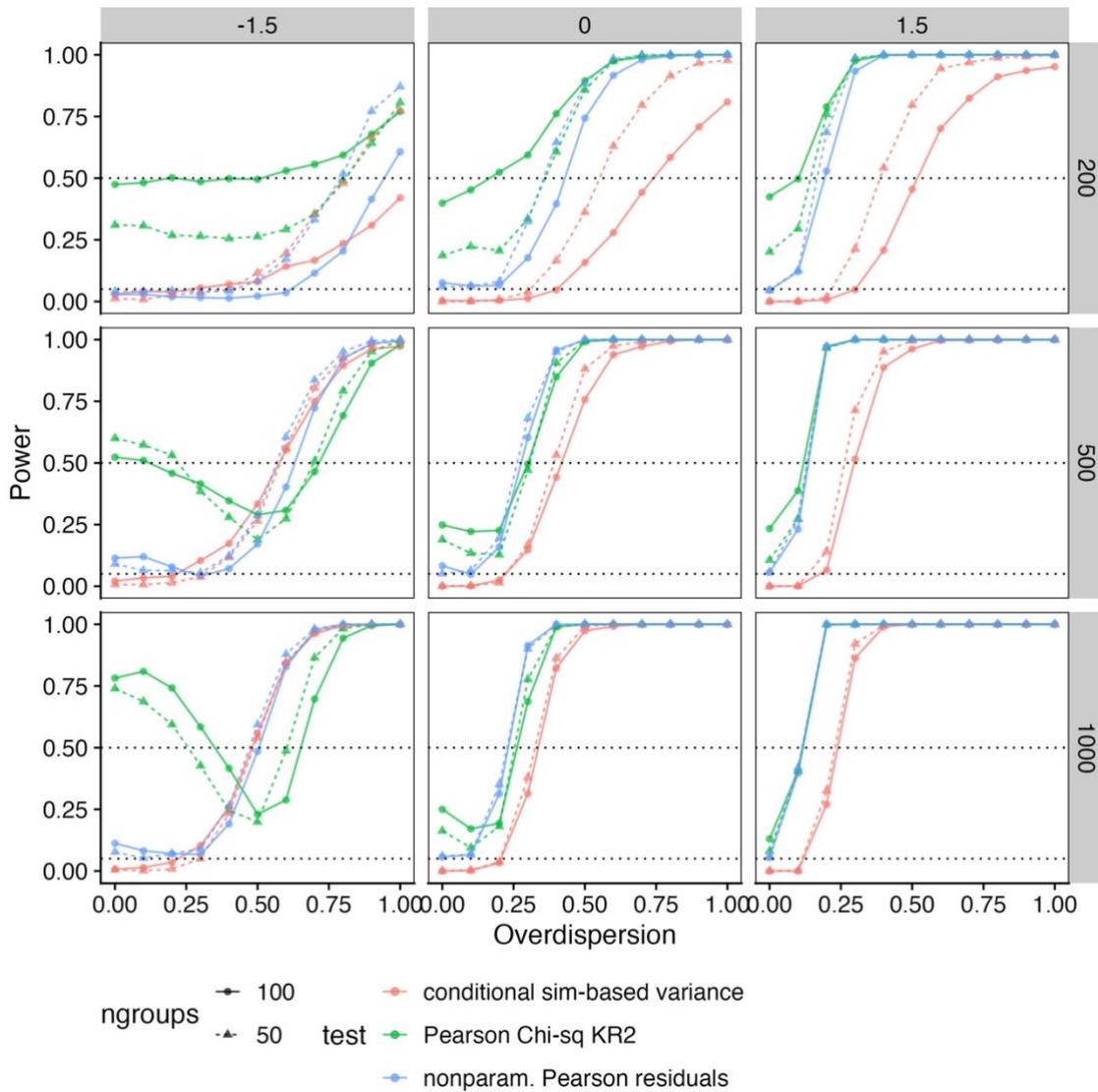
349

350 **Figures S8.2.** Type I error for the parametric Pearson test for Poisson GLMMs  
 351 performed with different corrections for the residual degrees of freedom (panel  
 352 columns), number of groups in the random intercept (panel rows) and sample size (x-  
 353 axis). Data were simulated from a Poisson GLMM with different intercepts (colours).  
 354 Please refer to the main text above to relate to each applied correction. 1000 simulations  
 355 for each parameter setting, slope = 1, random intercept variance = 1.



356

357 **Figure S8.3.** Dispersion parameters for the parametric Pearson test for Poisson GLMMs  
 358 performed with different corrections for the residual degrees of freedom (colours),  
 359 number of groups in the random intercept (linetype and shape), sample size (panel  
 360 rows), and intercept (panel columns). Please refer to the main text above to relate to  
 361 each applied correction. To improve clarity, we omitted the other corrections because  
 362 they are too similar to each other. 1000 simulations for each parameter setting, slope =  
 363 1, random intercept variance = 1.



364

365 **Figure S8.4.** Power of dispersion tests for Poisson GLMMs (colours) performed with  
 366 different numbers of groups in the random intercept (linetype and shape), sample size  
 367 (panel rows), and intercept (panel columns). Please refer to the main text above to relate  
 368 to the applied correction for residual degrees of freedom. To improve clarity, we omitted  
 369 other corrections for residual degrees of freedom because they are too similar to each  
 370 other. 1000 simulations for each parameter setting, slope = 1, random intercept variance  
 371 = 1.