3

5

6

8

9

10 11

12

13

14

15

16

17

18

19 20



Robustness of Statistical Power in Group-Randomized Studies of Mediation Under an Optimal Sampling Framework

Kyle Coxo and Benjamin Kelcey

Quantitative and Mixed Methods Research Methodologies, University of Cincinnati, OH, USA

Abstract: When planning group-randomized studies probing mediation, effective and efficient sample allocation is governed by several parameters including treatment-mediator and mediator-outcome path coefficients and the mediator and outcome intraclass correlation coefficients. In the design stage, these parameters are typically approximated using information from prior research and these approximations are likely to deviate from the true values eventually realized in the study. This study investigates the robustness of statistical power under an optimal sampling framework to misspecified parameter values in group-randomized designs with group- or individual-level mediators. The results suggest that estimates of statistical power are robust to misspecified parameter values across a variety of conditions and tests. Relative power remained above 90% in most conditions when the incorrect parameter value ranged between 50% and 150% of the true parameter.

Keywords: optimal sample allocation, power, multilevel mediation, group-randomized study

When investigating a treatment that operates within a hier- $\frac{22}{21}$ archical structure, experimental designs that assign groups 24 to treatment conditions often ameliorate ethical concerns, 25 more appropriately suit extant organizational structures, 26 and better reflect the contexts in which the results would 27 be generalized while maintaining a rigorous basis for causal inference (Raudenbush, 1997; Spybrook & Raudenbush, 28 2009). One of the primary questions in the planning stages 29 of these studies is the sample size required to detect the 30 31 expected treatment effect with a reasonably high probabil-32 ity. Sampling more groups almost always increases the power or probability to detect an effect in group-rando-33 34 mized studies (see Figure 1A) but this strategy typically 35 incurs substantial additional costs (Kelcey & Phelps, 2013, 36 2014; Raudenbush, 1997). Optimal sample allocation 37 frameworks balance these concerns by introducing cost parameters in the design phase that incorporate cost con-38 39 siderations directly into study planning. In turn, such frameworks identify a sample of individuals per group and groups 40 41 that maximize power given study-specific conditions and a 42 designated budget and cost structure (see Figure 1B; e.g., 43 Hedges & Borenstein, 2014; Kelcey, Phelps, Spybrook, 44 Jones, & Zhang, 2017).

Despite the conceptual utility of the optimal sampling 45 framework, a key constraint in its practical implementation 46 is a priori knowledge of the parameter values that govern 47 power for a particular estimand. In the planning phase, 48 the values of these parameters are typically not precisely 49 50 available, so researchers may utilize previous empirical results to infer a range of plausible values. It is likely that 51 these initial parameter estimates will deviate to some extent 52 53 from the true parameter values because estimates from the literature are almost always based on outcomes and popu-54 lations that differ in overt and subtle ways from the out-55 come and sample to be used in the new study (e.g., 56 57 Korendijk, Moerbeek, & Maas, 2010).

58 For these reasons, prior literature has investigated the robustness of designs in terms of their relative efficiency 59 60 to various types of parameter value misspecifications in the planning phase (e.g., Korendijk et al., 2010). This 61 research has largely been confined to the loss of efficiency 62 regarding estimation of the main effect. Recent research 63 has, however, promoted studies that are intentionally 64 designed to probe a more comprehensive set of effects that 65 unpack the theory of action underlying the treatment 66 effects using mediation (e.g., Gottfredson et al., 2015; 67



Figure 1. Power of a group-randomized study to detect group-level mediation using the Sobel test as a function of (A) group sample size (n_2) and (B) individual sample size (n_1) under an optimal sample allocation framework. *True optimal individual sample size while a dot marks the optimal individual sample size based on a misspecified value.

Institute of Education Sciences, US Department of Education, & National Science Foundation, 2013; Authors, B
[Author: Please provide full reference details]). Yet, little
is known regarding the robustness of power in these
designs when parameter values are misspecified.

73 In this study, we advance the literature base regarding the robustness of multilevel designs to parameter value 74 75 misspecifications in the context of mediation and optimal 76 sampling. More specifically, we incorporate study costs 77 and budget into the planning phase using an optimal sam-78 pling framework in order to contrast the power of two com-79 peting designs - one design that yields the maximum level 80 of power and an alternative design that yields a sub-optimal 81 level of power because of parameter value misspecifica-82 tions. We limit our considerations regarding sampling 83 design to the relative sample sizes at the individual- and 84 group-level.

85 To illustrate these differences, consider Figure 1B, it presents power as a function of the individual-level sample size 86 87 (n_1) under an optimal sampling framework. The peak of the power curve represents the most efficient use of resources 88 89 (i.e., sample of individuals per group and sample of groups) 90 which maximizes attainable power given study-specific conditions. An "*" identifies this optimal design and serves as 91 92 an appropriate benchmark for comparing alternative 93 designs because any deviation from the optimal design is 94 an inefficient use of limited resources and achieves less 95 power. This study tracks the consequences of using mis-96 specified parameter values in the study planning phase by 97 comparing power rates for an optimally designed study (* 98 in Figure 1B) to one planned with misspecified parameter 99 values (• in Figure 1B). Consequences are represented by 100 the distance between the true (*) and misspecified (•) opti-101 mal individual sample allocation in Figure 1.

102 Specifically, we estimate the power of a study design 103 when true parameter values were used to determine the optimal sample allocation and compare it to power rates 104 when the study employs a sample allocation based on mis-105 specified parameter values. We probe studies examining 106 the effect of a group-level treatment on an individual-level 107 outcome through a group-level mediator or an individual-108 level mediator with different sampling cost structures and 109 three tests of the mediation effect. We refer to these medi-110 ated effects by the level of the treatment, mediator, and 111 outcome (i.e., 2-2-1 and 2-1-1 mediated effects). Sampling 112 design considerations are limited to the relative sample 113 sizes at the group- and individual-level and assume an 114 equal number of groups in each of the two treatment con-115 ditions. The paper is divided into two major sections based 116 on these design types (i.e., 2-2-1 and 2-1-1). Both represent 117 investigations of multilevel mediation as the mediated 118 effect operates across levels of a hierarchy with the numer-119 ical notation identifying the level of the treatment-media-120 tor-outcome variables, respectively. For each design, we 121 122 outline the key components of the optimal sampling framework, describe the scope of our investigation, and detail our 123 results. We conclude with a brief discussion. 124

Group-Level Mediator

Analytic Model

In 2-2-1 designs groups are randomized to one of two treat-127 128 ment conditions (T) to examine treatment impacts on an individual-level outcome (Y) operating through a group-129 level mediator (M). We adopt a typical multilevel mediation 130 path formulation with parameters estimated using maxi-131 mum likelihood (e.g., Authors, D[Author: Please provide 132 full reference details]; Zhang, Zyphur, & Preacher, 133 2009). The mediator path model at the group-level can 134 be expressed as 135

125

126

Mediator model (Level 2) : $M_i = \pi_0 + aT_i + \varepsilon_i \varepsilon_i$

137
$$\sim N(0, \sigma_{M|}^2).$$
 (1)

138 The mediation model (i.e., Equation 1) captures the treat-139 ment-mediator relationship and describes variation in the 140 group-level mediator as a function of a group-level treat-141 ment. We use M_i as the mediator for group *j*, T_i as the treat-142 ment assignment coded as 0.5 and -0.5 for the treatment and control condition, respectively. Coefficient a represents 143 144 the treatment-mediator relationship with e_i as the normally distributed mediator error with a mean of zero and variance 145 σ_{M}^2 conditional on T_i . 146

147 The outcome model captures the relationship between148 the group-level mediator and individual-level outcome such149 that

Outcome model (Level 1) :
$$Y_{ij} = \beta_{0j} + e_{ij} e_{ij}$$

 $\sim N(0, \sigma_Y^2),$ (2a)

152

154

174

151

(Level 2):
$$\beta_{0j} = \gamma_{00} + bM_j + c T_j + u_{0j} u_{0j}$$

 $\sim N(0, \tau_{Y|}^2).$ (2b)

Here, Y_{ii} is the outcome for individual *i* in group *j*. At the 155 individual-level, e_{ij} represents the normally distributed error 156 157 term with a mean of zero and variance σ_v^2 . At the group-158 level, b is the conditional relationship between the mediator 159 and the outcome, while c' captures the direct effect of the 160 treatment on the outcome while controlling for the mediator, and u_{0i} is the normally distributed group-specific ran-161 dom effect with a mean of zero and variance τ_{V}^2 162 conditional on T_i and M_i . A combined mediator and out-163 164 come model from Equations 1 and 2can be expressed as

166
$$Y_{ij} = (\gamma_{00} + b\pi_0) + (ba + c')T_j + b\varepsilon_j + u_{0j} + e_{ij}.$$
 (3)

167 Finally, we utilize a standardized outcome and mediator 168 such that each has a mean of zero and unit variance. Under 169 this formulation $\sigma_M^2 = 1$ and $\tau_Y^2 + \sigma_Y^2 = 1$. This implies 170 $\tau_Y^2 = \rho$ and $\tau_Y^2 + \sigma_Y^2 = \rho + (1 - \rho)$ when ρ is the uncondi-171 tional intraclass correlation coefficient for the outcome 172 defined as

 $\rho = \frac{\tau_Y^2}{\sigma_Y^2 + \tau_Y^2}.$ (4)

175 Conditional and unconditional ρ values capture variance
176 in the outcome attributable to the group-level or the corre177 lation among individuals on the outcome within the same
178 group.

179This formulation places the *a* and *c'* path coefficients on a180Cohen's *d* or standardized group differences scale and

formats the b path coefficient as a standardized regression181coefficient.182

Mediation Test Statistics and Power

We estimate the 2-2-1 multilevel mediation effect (ME) 184 using the typical product of coefficients method: ME = ab185 with the *a* coefficient representing the treatment-mediator 186 path (Equation 1 and the b coefficient representing the 187 mediator-outcome path (Equation 2b. To determine the sig-188 nificance of this effect we consider the Sobel, joint, and 189 Monte-Carlo interval tests because they can be employed 190 before the collection of data (i.e., during study planning). 191 Performance of these tests in terms of power converges 192 with large sample sizes but under sample sizes typical for 193 studies in the social sciences there are substantial differ-194 ences (e.g., Authors D[Author: Please provide full refer-195 ence details]). Additionally, the optimal sample allocation 196 for a study can differ based solely on the selected mediation 197 test (Authors A[Author: Please provide full reference 198 199 details). These factors necessitate the inclusion of multiple mediation tests representing various approaches. 200

The historically popular Sobel Test compares the ratio of 201 the estimated mediation effect to its estimated standard 202 error. In recent years, the test has been heavily criticized 203 because of its imprecision in small-to-moderate sample 204 sizes. We outline the test for illustrative and comparative 205 purposes and recommend the use of alternative tests subsequently outlined. The Sobel test statistic is 207

$$z_{ab}^{\text{Sobel}} = ab/\sqrt{\sigma_{ab}^2}.$$
 (5) 209

Here, σ_{ab}^2 represents the error variance of the mediated effect and ab is the ME defined above. Based on prior literature we can estimate this error variance as a function of individual paths and their individual error variances and form the Sobel Test statistic for the mediation effect as (Sobel, 1982) 210

$$z_{ab}^{\text{Sobel}} = ab/\sqrt{b^2\sigma_a^2 + a^2\sigma_b^2}.$$
 (6) 217

The test statistic has an asymptotically normal distribution allowing inferences with large sample sizes to be drawn based on a comparison of the test statistic (e.g., z_{ab}^{Sobel}) to a critical value (e.g., $z_{\text{critical}} = 1.96$) in a standard normal distribution (Φ) at the associated type one error rate (e.g., 0.05 for 1.96).

When the alternative hypothesis is true, the test statistics224follow a non-central distribution with the ratios as the non-
centrality parameter. Power to detect the mediation effect225is then determined with227

183

$$P(|z_{ab}^{\text{Sobel}}| > z_{\text{critical}}) = 1 - \Phi(z_{\text{critical}} - z_{ab}^{\text{Sobel}}) + \Phi(-z_{\text{critical}} - z_{ab}^{\text{Sobel}}).$$
(7)

An alternative to the Sobel test is the joint test which 230 231 avoids the direct estimation of the mediation effect and 232 its distribution. Rather, it uses a composite null approach 233 that considers the treatment-mediator and mediator-out-234 come paths separately to determine the significance of a 235 mediated effect. The two concurrent sub-tests compare 236 the ratio of the path estimate and its standard error to a 237 normal or t-distribution. Only after rejecting the null 238 hypotheses of each individual path is the mediation effect 239 considered statistically significant. Avoiding the distributional assumptions of the Sobel test allows the joint test 240 to perform well in terms of power and type one error rate 241 242 in a variety of settings (e.g., Hayes & Scharkow, 2013; Kel-243 cey, Dong, Spybrook, & Shen, 2017).

The statistical test of the *a* path representing the treatment-mediator association is

$$z_a = a/\sigma_a, \tag{8}$$

and the test of the *b* path representing the mediator-out-come association is

$$z_b = b/\sigma_b. \tag{9}$$

252 With σ_a^2 and σ_b^2 indicating the error variances of the *a* 253 path and *b* path respectively. Power to detect the mediation 254 effect using the joint test is simply the product of the power 255 to detect each path which we formulate as

$$P(|z_a| > z_{\text{critical}} \text{ and } |z_b| > z_{\text{critical}})$$

$$= (1 - \Phi(z_{\text{critical}} - z_a) + \Phi(-z_{\text{critical}} - z_a))$$

$$\times (1 - \Phi(z_{\text{critical}} - z_b) + \Phi(-z_{\text{critical}} - z_b)), \quad (10)$$

258 with Φ () as the normal cumulative density function.

259 The Monte-Carlo (MC) interval test resamples a and b260 path values with sampling variability equal to the error vari-261 ance of the respective paths (Preacher & Selig, 2012). Max-262 imum likelihood estimation identifies path coefficient estimates with an assumed multivariate normal distribution 263 264 with means, variances, and covariances based on the max-265 imum likelihood estimates (Preacher & Selig, 2012). For the 266 2-2-1 analytic models described above, we draw a and b 267 path values using

269
$$\begin{pmatrix} a^{\#} \\ b^{\#} \end{pmatrix} \sim \text{MVN}\left(\begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}, \begin{pmatrix} \hat{\sigma}_{\hat{a}}^2 & \hat{\sigma}_{\hat{a},\hat{b}} \\ \hat{\sigma}_{\hat{a},\hat{b}} & \hat{\sigma}_{\hat{b}}^2 \end{pmatrix} \right).$$
(11)

270 An approximation of the sampling distribution of the 271 mediated effect is formed using the product of each set 272 of $a^{\#}$ and $b^{\#}$. Inferences regarding the mediated effect are 273 drawn based on the inclusion of zero in the asymmetric confidence intervals constructed from the sampling distri-
bution. Power of the MC interval test is simply the propor-
tion of asymmetric confidence intervals that exclude zero.275276

Optimal Sample Allocation

Theoretically, optimal sample allocation provides research-278 ers a means to identify the sampling strategy that maxi-279 mizes power under specific constraints on design and 280 budget (e.g., Kelcey, Phelps, et al., 2017). The process 281 begins with the identification of the optimal sample of indi-282 viduals per group (n_{i}^{opt}) and then a simple function of this 283 value, budget, and sampling costs identifies the optimal 284 number of groups n_2^{opt} . 285

Throughout this study, we apply the conventional linear 286 cost formulation (Raudenbush, 1997) such that 287

$$T = c_2 n_2 + c_1 n_2 n_1. \tag{12}$$

277

where T is the total funds available to collect data for a 290 study, c_1 is the cost to enroll each individual after sam-291 pling the group, and c_2 is sampling cost for each additional 292 group. Each is typically measured in monetary units (e.g., 293 dollars). For our optimal sample allocation formulas, we 294 assume an equal number of groups in two treatment con-295 ditions (i.e., treatment and control conditions) and equal 296 cost across conditions. 297

Like group-randomized studies of main effects, optimal298sample allocation under the Sobel test is derived by mini-
mizing the error variance (Raudenbush, 1997). With a stan-
dardized mediator and outcome, the optimal individual
sample size for the Sobel test in terms of path coefficients300sample size for the Sobel test in terms of path coefficients302is (Authors, A[Author: Please provide full reference303details])304

$$n_{1}^{\text{opt}} = \sqrt{\binom{c_{2}}{c_{1}}} \frac{a^{2}(1-\rho)}{4b^{2}\left(1-\frac{a^{2}}{4}\right)^{2} + a^{2}\left(\rho - \frac{\left(ab+c'\right)^{2}}{4} - b^{2}\left(1-\frac{a^{2}}{4}\right)\right)},$$
(13) 306

with the optimal sample of groups determined by substituting n_1^{opt} in $n_2^{\text{opt}} = T/(c_2 + c_1 n_1)$. 308

The structure of the joint test complicates the calculation 309 of an optimal individual sample size because there is no 310 mediation effect error variance to minimize. Rather, one 311 must directly maximize power to determine n_1^{opt} under 312 the joint test. There is no simple closed-form expression, 313 but one is possible through numerical methods (see Elec-314 tronic Supplemental Material, ESM 1; Authors, A [Author: 315 Please provide full reference details]). 316

Like the joint test, the MC interval test does not have a 317 closed form solution for n_1^{opt} but can be approximated by 318 the Sobel and joint test formulations in larger samples 319

320 (Authors, A). However, there are typical conditions (e.g., 321 large sampling cost ratios; c_2/c_1 in which the n_1^{opt} values under the Sobel and joint test substantially diverge from 322 323 the values for the MC interval test (Authors, A). Such con-324 ditions would introduce a confounding effect into this 325 investigation because we could not determine if differences 326 in power from inefficient sample allocation were due to 327 parameter value misspecification or a misalignment 328 between the Sobel and joint test n_1^{opt} and the true n_2^{opt} for 329 the MC interval test. We avoid these confounding effects 330 by numerically estimating n_1^{opt} for the MC interval test using 331 a linear search algorithm under the specified design and 332 budget.

Power of 2-2-1 Mediation Studies

334 Path coefficients (a, b, c'), the unconditional intraclass cor-335 relation coefficient (ρ), and cost structure (c_2/c_1) all influ-336 ence optimal sample allocation. Given that the parameter values employed in the study design phase will not precisely 337 338 match their true values, we compare power rates when opti-339 mal sample allocation is determined using true and mis-340 specified parameter values. Power is determined 341 analytically using the tests and formulations above and path 342 coefficient error variance formulations that utilize budget 343 and sampling cost values. Assuming balanced random assignment, we formulate these error variances such that 344 345 (Authors, A; Authors D [Author: Please provide full reference details]) 346

349 and

348

351

$$\sigma_b^2 = \frac{\left(\rho - \frac{(ab+c')^2}{4} - b^2\left(1 - \frac{a^2}{4}\right)\right) + (1 - \rho)/n_1}{\left(1 - \frac{a^2}{4}\right)(T/(c_2 + c_1n_1))}.$$
 (15)

 $\sigma_a^2 = \frac{4(1 - (a^2 + 4))}{T/(c_2 + c_1 n_1)}$

These formulations employ parameters common to planning group-randomized mediation studies (Authors, B; Authors D) and, under an optimal sampling framework, allow power analyses using only n_1^{opt} . Cost and budget information has replaced n_2 values.

357 Borrowing from Korendijk et al. (2010), we refer to the true but unknown parameter values as population values, 358 denoting them with * (e.g., a^* , b^* , c'^* , ρ^*), and labeling 359 360 the true optimal individual per group sample size as $n_1^{\text{opt}*}$. The estimated or predicted misspecified values used in 361 362 the study design phase are referred to as the initial values 363 (e.g., a, b, c', ρ) with the optimal individual per group sam-364 ple size based on these values retaining the n_1^{opt} notation.

We defined the robustness of statistical power under an
optimal sampling framework against misspecified parame-
ter values in terms of the relative loss of statistical power:365
366367

Relative Power =
$$\frac{\text{Initial Design Study Power}}{\text{Population Design Study Power}}.$$
(16) 369

To determine the power of a design under the population 370 values, we identified the optimal individual sample size 371 using the true population values $(n_1^{\text{opt}*})$ and then calculated 372 the implied power to detect the mediation effect. Similarly, 373 to determine the power of a design under the initial values, 374 we identified the optimal individual sample size using the 375 initial misspecified parameter values (n_1^{opt}) and then calcu-376 lated the implied power to detect the mediation effect. Rel-377 ative power examines the loss of power associated with an 378 n_1^{opt} based on a misspecified parameter value and has a 379 maximum value of one when the initial and population 380 parameters are equal (i.e., $n_1^{\text{opt}*} = n_1^{\text{opt}}$). Imprecision in our 381 linear search algorithm for the MC interval test did produce 382 some relative power values exceeding one but these erro-383 neous values only reflect noise in the estimation of initial 384 and population MC interval test power. They do not bias 385 the overall robustness of power results or interpretations. 386

Without clear benchmarks regarding what constitutes a 387 scale for relative power, we borrowed from the optimal 388 design literature regarding relative efficiency and interpret 389 a relative power value of .9 and above (i.e., 10% loss of 390 power or less) to be good and values between .8 and .9 391 as acceptable (e.g., Korendijk et al., 2010). Our general 392 descriptions of robustness are based on relative power 393 meeting or exceeding these benchmarks when the misspec-394 ified parameter values are between 50% and 150% of the 395 true population value. Specific levels of acceptable power 396 loss and parameter misspecification are a study-specific 397 consideration. 398

Our investigation included three different path coeffi-399 cient ratios (a/b = 0.3/0.3, 0.5/0.3, 0.6/0.2), three intra-400 class correlation coefficient values ($\rho = .3, .4, .5$), and 401 three group to individual sampling cost structures ($c_2/c_1 =$ 402 5/1, 10/1, 100/1). We employed a fully crossed design 403 resulting in a total of 27 conditions with various power rates 404 and n_1^{opt} values. Under each condition, we found relative 405 power when incorrect initial a values ranged from 0.05 to 406 0.9, incorrect initial b values ranged from 0.05 to 0.5, 407 and incorrect initial intraclass correlation coefficient values 408 ranged from 0.15 to 0.9 all using an interval of 0.025. 409 While not comprehensive, these explicit comparisons 410 reflect a range of substantive applications. It is also possible 411 to examine the robustness of statistical power to optimal 412 sample allocation based on misspecified parameters values 413 analytically. We found the analytic expressions to be 414 415 complex and unable to provide clear, intuitive, and descrip416 tive context for the properties of statistical power robust417 ness in these settings when compared to our explicit
418 comparisons. We provide an example analytic analysis in
419 ESM 1.

420 While the direct effect or c' is also included in the optimal 421 sample allocation formulas it has little influence on n_1^{opt} and 422 relative power. It is held constant at c' = 0.1 throughout the 423 analyses and not discussed further (see ESM 1). Addition-424 ally, the total funds (*T*) do not play a role in determining 425 the optimal individual sample size so *T* was set as a func-426 tion of group cost (*T* = 100 c_2 ; see Authors, A).

427 **Results**

428 Results are presented by misspecified parameter (*a*, *b*, and 429 ρ) with a focus on the general patterns involving robustness 430 of power to capture the essences of the analysis. The com-431 plete results for the fully crossed design can be found in 432 ESM 1, Tables 2, 3, and 4.

433 Incorrect Initial Intraclass Correlation Coefficient

We examined power when n_{i}^{opt} was determined using ρ val-434 435 ues ranging from .15 to .9 with true population values of ρ^* 436 = .3, .4, and .5. Overall, power was very robust to incorrect ρ 437 values under stated conditions. With ρ from 50% to 150% 438 of ρ^* , relative power remained greater than .9 in nearly every condition (see Table 1 for an example of complete 439 440 results when $\rho^* = .5$). This indicates designs with a misspec-441 ified ρ value suffered less than a 10% loss in power due to 442 inefficient sample allocation.

A few noteworthy patterns emerged from these results. 443 444 First, larger a/b path coefficient ratios decreased the 445 robustness of power to incorrect p values. Second, over-446 and underestimated p values had similar influence on relative power. In other words, the loss of power associated 447 448 with using incorrect ρ values to determine n_1^{opt} was fairly 449 symmetric. This symmetric power loss was consistent 450 across each test. Lastly, when the group-to-individual cost ratio (c_2/c_1) was small, power was more susceptible to inef-451 ficient sample allocation from misspecified p values (see 452 Figure 2B). As c_2/c_1 increased, designs became more robust 453 454 with relative power remaining well above .9 across a wide 455 range of predicted p values in designs with the highest cost 456 ratio.

457 Incorrect Initial *a* Path Coefficient

458 Our examination found relative power when n_1^{opt} was deter-459 mined using *a* values ranging from .05 to .9 with true pop-460 ulation values of $a^* = .3$, .5, and .6. Overall, power was very 461 robust to incorrect initial *a* path coefficient values under 462 stated conditions. When *a* ranged from 50% to over 463 150% of a^* , relative power remained greater than .9 in almost every condition. While results are similar, power 464 was generally more robust when misspecified *a* values were 465 used to determine n_1^{opt} and less robust using misspecified ρ 466 values. 467

Effects of misspecified a values (see Figure 3B) varied by 468 mediation test with designs using the MC interval test pro-469 viding the greatest degree of robustness followed closely by 470 the joint test. Relative power in designs utilizing these tests 471 remained above .95 across nearly the full range of *a* values. 472 One exception were the reductions in relative power under 473 the joint test when a overestimated a^* . The inverse was 474 true when using the Sobel test, inefficient sample allocation 475 caused by a values that underestimated a^* were more detri-476 mental to power. 477

Similar to our examination of incorrect ρ values, power was more robust to incorrect *a* values when the design had larger (c_2/c_1) values (see Figure 3B). Additionally, power in designs with greater ρ values was more robust to inefficient sample allocation based on incorrect *a* values. 482

Incorrect Initial b Path Coefficient

We examined power when n_1^{opt} was determined with b values from .05 to .5 and true population values of $b^* = .2$ and484.3. Like previous results, power under the optimal sampling486framework was robust to incorrect b values. When b was487anywhere from 50% to 150% of b^* relative power remained488above .9.489

483

502

503

Overall, c_2/c_1 and ρ values did little to influence the 490 robustness of power when incorrect b values were used to 491 determine n_1^{opt} but we again found relative power varied 492 by test (see Figure 4B). Relative power using the Sobel test 493 and MC interval test was similar with notable differences 494 occurring mostly in designs with smaller p values and larger 495 c_2/c_1 values. Under those conditions power using the Sobel 496 test decreased at a greater rate when the b value used to 497 determine n_1^{opt} underestimated b^* . Designs using the joint 498 test also demonstrated robustness to incorrect b values 499 but when b significantly overestimated b^* inefficient sample 500 501 allocation caused power rates to fall sharply.

Individual-level Mediator

Analytic Model

We next consider a group-randomized design with an indi-504 vidual-level mediator. Here, we again have intact groups 505 assigned to a treatment condition (T) and an individual-506 level outcome (Y) but the relationship under examination 507 flows through an individual-level mediator (M). For this 508 design parameters are also estimated using maximum like-509 lihood but a multilevel model is needed for the media-510 tor such that (Authors, B; Pituch & Stapleton, 2012; 511

Conditio	ns		ρ									
C ₂ :C ₁	a:b	Test	.2	.3	.4	.5*	.6	.7	.8	.9		
		S	.993	.997	.999	1.000	.999	.993	.978	.935		
	0.3:0.3	JT	.986	.992	.998	1.000	.996	.981	.943	.840		
		MC	.976	.955	1.006	.980	1.000	1.012	.918	.828		
		S	.986	.993	.998	1.000	.997	.988	.963	.898		
5	0.5:0.3	JT	.986	.992	.998	1.000	.996	.982	.945	.848		
		MC	.993	.995	1.011	1.006	.994	1.005	.952	.848		
		S	.966	.985	.996	1.000	.995	.979	.942	.854		
	0.6:0.2	JT	.937	.973	.993	1.000	.993	.969	.919	.809		
		MC	.923	.944	.985	1.004	.992	.984	.943	.786		
		S	.995	.997	.999	1.000	.999	.994	.983	.948		
	0.3:0.3	JT	.996	.997	.999	1.000	.998	.989	.963	.885		
		MC	1.025	1.014	1.042	.997	1.031	1.015	.987	.879		
		S	.989	.995	.999	1.000	.998	.990	.971	.918		
10	0.5:0.3	JT	.997	.998	.999	1.000	.998	.990	.965	.893		
		MC	1.003	1.009	1.002	1.012	.991	.985	.983	.876		
		S	.972	.988	.997	1.000	.996	.983	.952	.876		
	0.6:0.2	JT	.953	.979	.995	1.000	.994	.974	.930	.830		
		MC	.955	.993	1.006	1.016	1.014	.997	.974	.839		
		S	.998	.999	1.000	1.000	1.000	.998	.993	.980		
	0.3:0.3	JT	.997	.999	1.000	1.000	1.000	1.000	.997	.984		
		MC	1.002	1.031	1.011	1.051	1.017	1.023	1.023	1.006		
		S	.996	.998	.999	1.000	.999	.996	.989	.969		
100	0.5:0.3	JT	.996	.999	1.000	1.000	1.000	1.000	.997	.985		
		MC	1.012	1.002	.999	1.003	1.018	1.008	1.003	.991		
		S	.989	.995	.999	1.000	.998	.993	.980	.945		
	0.6:0.2	JT	.996	.995	.998	1.000	.998	.990	.972	.926		
		MC	1.008	1.006	1.008	1.015	.998	1.017	.988	.949		

Table 1. Relative power of group-randomized studies of 2-2-1 mediation with misspecified intraclass correlation coefficients when p*=.5

Notes. S = Sobel; JT = Joint, MC = Monte-Carlo. *True population parameter value. Full results with an interval of 0.025 between incorrect initial values are available in ESM 1, Table 2.

 512
 Raudenbush & Bryk, 2002; VanderWeele, 2010; Zhang

 513
 et al., 2009)

$$M_{ij} = \pi_{0j} + \varepsilon^{M}_{ij} \ \varepsilon^{M}_{ij} \sim N(0, \sigma^{2}_{M}),$$

515
$$\pi_{0j} = \zeta_{00} + aT_{j} + u^{M}_{0j} \ u^{M}_{0j} \sim N(0, \tau^{2}_{M|}),$$
 (17)

516 with a combined mediator model formulation 517 expressed as

519
$$M_{ij} = \zeta_{00} + aT_j + u_{0j}^M + \varepsilon_{ij}^M, \qquad (18)$$

520 where M_{ij} represents the mediator value for individual *i* in group j, T_i as the treatment assignment coded as 0.5 and 521 522 -0.5 for the treatment and control condition, respectively 523 with associated path coefficient *a*, ε_{ij}^{M} as the normally distributed error term with a mean of zero and variance $\sigma_{M}^{\scriptscriptstyle 2},$ 524 and u_{oi}^{M} as the group-specific random effects that follow a 525 526 normal distribution with mean zero and variance $\tau_{M|}^2$ con-527 ditional on T_i . Conceptually, *a* maps out how exposure to the treatment produces changes in the individual-level 528 mediator. For the 2-1-1 design, the outcome model is 529

$$\begin{aligned} Y_{ij} &= \beta_{0j} + b_1 \big(M_{ij} - \bar{M}_j \big) + \epsilon_{ij}^Y \; \epsilon_{ij}^Y \sim N \Big(0, \sigma_{Y|}^2 \Big), \\ \beta_{0j} &= \gamma_{00} + B \bar{M}_j + c' T_j + u_{0j}^Y \; u_{0j}^Y \sim N \Big(0, \tau_{Y|}^2 \Big), \end{aligned} \tag{19}$$

with the combined outcome model formulation expressed 532 as 533

$$Y_{ij} = \left(\gamma_{00} + B\bar{M}_j + c'T_j + u_{0j}^{Y}\right) + b_1\left(M_{ij} - \bar{M}_j\right) + \varepsilon_{ij}^{Y},$$

$$Y_{ij} = \gamma_{00} + (B - b_1)\bar{M}_j + b_1M_{ij} + c'T_j + u_{0j}^{Y} + \varepsilon_{ij}^{Y}.$$
(20) 535

We retain Y_{ij} as the outcome for individual *i* in group *j*, 536 and use $M_{ij} - \bar{M}_j$ as the group-centered individual-level 537 mediator with coefficient b_1 , \bar{M}_j as the mean of the mediator in group *j* with path coefficient *B*, *c'* as the treatmentoutcome conditional path coefficient, and u_{oj}^Y and ε_{ij}^Y as 540



Figure 2. Relative power of group-randomized studies of 2-1-1 (A) and 2-2-1 (B) mediation using the Sobel (dash), joint (dot), and MC interval test (solid) as a function of incorrect intraclass correlation coefficient values with different group to individual cost ratios and a total budget of 100 times the cost ratio when a = .5, b (or B) = .3, c' = 1, and the true intraclass correlation value is a .4. Vertical lines mark incorrect initial parameter values that are 50% and 150% of the true value and a horizontal line marks the 90% relative power benchmark.



Figure 3. Relative power of group-randomized studies of 2-1-1 (A) and 2-2-1 (B) mediation using the Sobel (dash), joint (dot) and MC interval test (solid) as a function of incorrect a path coefficient values with different group to individual cost ratios and a total budget of 100 times the cost ratio when b (or B) = .3, c' = 1, all intraclass correlation values are .5, and the true a path coefficient value is a^* = .5. Vertical lines mark incorrect initial parameter values that are 50% and 150% of the true value (a^*) and a horizontal line marks the 90% relative power benchmark.



Figure 4. Relative power of group-randomized studies of 2-1-1 (A) and 2-2-1 (B) mediation using the Sobel (dash), joint (dot) and MC interval test (solid) as a function of incorrect *b* (or *B*) path coefficient values with different group to individual cost ratios and a total budget of 100 times the cost ratio when a = .5, c' = 1, all intraclass correlation values are .5, and the true b (or B) path coefficient value is 0.3. Vertical lines mark incorrect initial parameter values that are 50% and 150% of the true value (*b** or *B**) and a horizontal line marks the 90% relative power benchmark.

(22)

541 the level-two and level-one error terms. Both error terms 542 are normally distributed with variances of σ_{V}^2 conditional on $(M_{ii} - \bar{M}_i)$ for ε_{ii}^Y and $\tau_{Y_i}^2$ conditional on \bar{M}_i and T_i for 543 544 u_{oi}^{Y} . Conceptually, the mediator-outcome relationship occurring at the individual-level is represented by b_1 and 545 at the group-level by $(B - b_1)$ with B representing the total 546 547 individual- and group-level relationship. A combined medi-548 ator and outcome model for 2-1-1 mediation can be 549 expressed as

$$Y_{ij} = (\gamma_{00} + b_1 \zeta_{00}) + (B - b_1) \overline{M}_j + (b_1 a + c) T_j + b_1 u_{0j}^M + b_1 \varepsilon_{ij}^M + u_{0j}^Y + \varepsilon_{ij}^Y,$$
(21)

552 which allows a clear comparison of the 2-2-1 and 2-1-1 553 analytic models.

554 In group-randomized studies of 2-1-1 mediation, there 555 are two intraclass correlation coefficients. The first, ρ_{y_1} cor-556 responds to the ρ described for group-randomized studies 557 of 2-2-1 mediation and captures the correlation among indi-558 viduals on the outcome within the same group. The second, 559 ρ_M , is introduced as a result of the multilevel nature of the mediator model in the 2-1-1 design. It captures the correla-560 561 tion among individuals on the mediator within the same group and we formulate the unconditional intraclass corre-562 563 lation coefficient of the mediator as

$$ho_M=rac{ au_M^2}{\sigma_M^2+ au_M^2}.$$

We use a standardized outcome and mediator similar to that described for ρ under the 2-2-1 model resulting in similar scaling for the *a*, *c'*, *b*₁, and *B* paths (i.e., *a* and *c'* as standardized group differences and *b*₁ and *B* as standardized regression coefficients). 560

Our utilization of group-mean centering operationalizes B 571 to capture the total influence of the mediator on the out-572 come. We focus on this cumulative or overall mediation 573 estimated using the typical product of coefficients method: 574 ME = aB (Pituch & Stapleton, 2012, VanderWeele, 2010; 575 VanderWeele & Vansteelandt, 2009). The 2-1-1 model also 576 captures the relationship between the treatment and out-577 come as it operates through a mediator at the group-578 $(B - b_1)$ and individual-level (b_1) . The $a(B - b_1)$ estimate 579 of the mediated effect in the 2-1-1 model is conceptually 580 equivalent to the ab mediated effect of the 2-2-1 model 581 (Pituch & Stapleton, 2012). 582

Power and Optimal Sample Allocation 583

The general form of the mediation tests and subsequent 584 power formulations presented for the 2-2-1 design are sim-585 ilar for 2-1-1 design. The error variances associated with 586 each path and the error variance of the mediated effects 587 do change substantially but formulations utilizing path coef-588 ficients, budget, and sampling cost are available in the liter-589 ature along with related optimal sample allocation 590 formulations (see ESM 1; Authors, B; Authors C). 591

551

565

592 Unlike 2-2-1 designs, there is no simple closed form solu-593 tion for n_1^{opt} under the Sobel test. We can, however, determine n_1^{opt} for the Sobel test numerically (see ESM 1; 594 595 Authors, C [Author: Please provide full reference 596 **details**]). Similar processes are used to determine n_1^{opt} 597 under the joint test for 2-1-1 and 2-2-1 designs. Using the 598 cost function formulations for n_1 and the error variance 599 path formulations, we directly maximize the power function (see ESM 1; Authors, C). The n_1^{opt} values from the Sobel and 600 601 joint test formulations approximate those of the MC inter-602 val test for 2-1-1 designs but we again use a linear search 603 algorithm to identify n_1^{opt} values for the MC interval test 604 to avoid possible confounding effects.

605 **Power of 2-1-1 Mediation Studies**

606 We now repeat our investigation but for group-randomized 607 studies of 2-1-1 mediation. For this design, we investigate 608 the robustness of power against misspecified a, B, ρ_M , 609 and ρ_{Y} values in the n_{y}^{opt} formulation. Procedures and con-610 ditions are similar to those described for the 2-2-1 design with the *B* parameter assuming the same values as those 611 612 assigned to b, and the unconditional intraclass correlation 613 coefficient of the mediator (ρ_M) and outcome (ρ_Y) assuming 614 the same values as those assigned to p. Results from an ini-615 tial simulation with differing values for the misspecified values were not qualitatively different from those presented 616 below so we constrained ρ_M and ρ_Y to be equal. That is, 617 $\rho_M^* = \rho_Y^*$ when we set the population values and $\rho_M =$ 618 619 ρ_{Y} as we varied the initial intraclass correlation coefficient 620 values. The initial simulation did suggest that the ρ_Y value 621 is likely the more influential parameter in terms of relative 622 power. Complete results for the 2-1-1 investigation can be 623 found in ESM 1, Tables 2, 3, and 4.

624 **Results**

625 Incorrect Initial Intraclass Correlation Coefficients

626 We examined power when n_1^{opt} was determined using ρ_M 627 and ρ_Y values ranging from .15 to .9 with true population 628 values of ρ_M^* and $\rho_Y^* = .3$, .4, and .5. Overall, power was 629 very robust to incorrect ρ_M and ρ_Y values under stated con-630 ditions. In nearly every situation, relative power remained 631 above .9 when ρ_M and ρ_Y ranged from 50% to nearly 632 200% of ρ_M^* and ρ_Y^* .

633 Over- or underestimating ρ_M^* and ρ_Y^* led to similar 634 amounts of lost power due to inefficient sample allocation 635 with power loss slightly greater in designs using the joint 636 test. Power was more robust to n_1^{opt} based on misspecified 637 parameter values when the design had higher sampling cost 638 ratios (i.e., c_2/c_1 ; see Figure 2A) or larger path coefficient ratios (i.e., a/b). The influence of the path coefficient ratio639was less pronounced when using the joint test. Additionally,640higher c_2/c_1 muted the benefits of larger path coefficient641ratios.642

Incorrect Initial a Path Coefficient

Our examination found relative power when n_1^{opt} was deter-644 mined using a values ranging from 0.05 to 0.9 with true 645 population values of $a^* = .3, .5, and .6$. Overall, power 646 was robust to misspecified a values that led to inefficiencies 647 in sample allocation. Across the different mediation tests 648 and conditions, a values ranging from 50% to 150% of a^* 649 almost always maintained relative power above .8. While 650 robust when compared to our 80% relative power bench-651 mark, misspecifying a values when determining the n_1^{opt} 652 for a group-randomized 2-1-1 mediation study led to the lar-653 gest reductions in power across the conditions and designs 654 considered here. 655

Results indicated that power decreased at a smaller rate 656 when a underestimated a^* in the determination of n_1^{opt} . 657 Conversely, power rates were less robust when a overesti-658 mated a^* when determining n_1^{opt} , especially in designs that 659 utilized the joint test (see Figure 3A). This relationship 660 became stronger in designs with larger ρ_M and ρ_Y values 661 and a/b ratios. Misspecified a values in the n_1^{opt} formulation 662 had similar detrimental effects on power across designs 663 with different a/b ratios but effects varied across designs 664 with different ρ_M and ρ_Y values. Here, relative power was 665 greater in designs with larger ρ_M and ρ_Y values. 666

For the *a* parameter, we found a unique relationship between relative power and c_2/c_1 . Power was more robust to misspecified *a* values in the n_1^{opt} formulation, when 2-1-1 designs had smaller c_2/c_1 ratios (see Figure 3A). In all other conditions and with different parameters, the converse was true. 672

Incorrect Initial B Path Coefficient

673

643

We examined power when n_1^{opt} was determined with *B* values ranging from 0.05 to 0.5 and true population values of674 $B^* = 0.2$ and 0.3. Following previous results, power under676the optimal sample allocation framework was very robust677to incorrect *B* values under stated conditions. Relative678power remained at or above .9 across tests and under679nearly every condition when *B* was 50–150% of *B**.680

When using the joint test or MC interval test, power was 681 even more robust to misspecified B values in the n_1^{opt} formu-682 lation and over a greater range of B values. However, in 683 designs that used the Sobel test, power decreased rapidly 684 due to inefficient sample allocation when B was less than 685 B^{*}. Larger c_2/c_1 , ρ_M and ρ_Y values did offset some of the 686 detrimental effects of incorrect B values on Sobel test 687 power. Designs using the joint or MC interval test had such 688 high relative power across the range of *B* values that their
parameter-power relationships were practically concealed
(see Figure 4A).

692 Influence of Mediation Test and Analytic Model

Results revealed several differences in the robustness of 693 694 power across analytic model, mediation test, and type of 695 misspecified parameter value which deserve some attention. First, the overall robustness of power under an optimal 696 697 sampling framework is less surprising when considering n_1^{opt} is typically small and relatively stable across these models 698 699 and the typical conditions included in this study (Author A, Author C). Put differently, n_1^{opt} is often small (e.g., 700 701 < 10) across a variety of typical conditions and true param-702 eter values. Therefore, using misspecified parameter val-703 ues, even with large discrepancies, still results in a small 704 n_1^{opt} values, likely < 10. Because the actual n_1^{opt} value 705 remains fairly stable, it is reasonable to find only minor 706 change in power rates. Additionally, previous research has 707 shown that study design conditions with less stable n_1^{opt} val-708 ues and therefore an increased likelihood of large differences between n_1^{opt} and $n_1^{\text{opt}*}$, are the same conditions in 709 which these differences have little influence on power 710 711 (Author A, Author C). In other words, when minor changes in parameter values influence n_1^{opt} values, deviations from 712 n_{1}^{opt} do not influence power and under conditions where 713 deviations from n_1^{opt} are detrimental to power, n_1^{opt} values 714 715 tend to be so similar (e.g., < 10) power is not substantially 716 influenced.

717 Another consistent result across models and tests was the influence of unconditional intraclass correlation coefficient 718 values (see Figure 2). The consistent influence of ρ (or ρ_M) 719 720 and $\rho_{\rm v}$) on the robustness of power under the optimal sampling framework is directly related to the consistent influ-721 722 ence of ρ on n_1^{opt} values (Author A, Author C). If the ρ 723 value is wrong (i.e., misspecified) then it is likely we have poorly estimated n_{i}^{opt} resulting in inefficiencies and a sub-724 725 stantial loss of statistical power. This is especially true in 726 designs with larger ρ values but less so in designs with large 727 c_2/c_1 ratios.

728 Relative power did vary across models and tests when 729 the path coefficients were misspecified with the exception 730 of the MC interval test for which power demonstrated con-731 sistent robustness. The MC interval test result can be traced back to stable n_1^{opt} values when using the test even with dif-732 733 ferent *a* and *b* path coefficient values (Author A, Author C). 734 Conversely, power under the optimal sampling framework 735 for the joint and Sobel test varied in robustness to misspec-736 ified a and b path coefficient values across models reflecting the conditions in which these parameters influenced 737 n_1^{opt} values for that specific test. For example, in a group-738 739 randomized study of 2-2-1 mediation using the joint test, 740 changes to the b path coefficient resulted in substantial 744

changes to n_1^{opt} values suggesting that misspecified b path741coefficient values will influence optimal sample allocation742and therefore power (Authors A).743

Illustrative Example

Pulling from an investigation of the Every Classroom, 745 Everyday (ECED) program by Early et al. (2016), we illus-746 trate the process described above. ECED aims to improve 747 student academic achievement, an individual-level out-748 come, through improvements in teaching practice and 749 curriculum alignment. The ECED program represents a 750 group-level treatment as it is implemented across whole 751 schools. For our illustration, we also include a group-level 752 mediator that captures the degree of program implementa-753 tion within a school, a crucial factor in the success of these 754 programs (Desimone, Porter, Garet, Yoon, & Birman, 755 2002). We now have a two-level group-randomized study 756 examining a 2-2-1 mediation effect. To test the significance 757 of this effect, we employ the joint test and for study plan-758 ning purposes predict the relationship between the ECED 759 program and our measure of program implementation is 760 a = 0.5, the relationship between program implementation 761 and student outcomes is b = 0.2, the conditional direct 762 effect of the ECED program on student outcomes is c' =763 0.1, the correlation of student outcomes within a school is 764 $\rho = 0.3$, with a budget of T = US\$175,000 and a school 765 to student sampling cost ratio of $c_2/c_1 = 2,000$. Under these 766 conditions, the optimal individual sample size is 38 students 767 per school with a sample of 86 schools (i.e., $n_1^{\text{opt}} = 38$ and 768 $n_2^{\text{opt}} = 86$). However, in the empirical results of our hypo-769 thetical study the true *a* path value was $a^* = 0.65$ indicating 770 we underestimated the *a* path during study planning and 771 therefore did not use the most efficient sample. The true 772 optimal design would have included a sample of 48 stu-773 dents per school in 85 schools (i.e., $n_1^{\text{opt}*} = 48$ and 774 $n_2^{\text{opt*}} = 85$). As conducted, power to detect the 2-2-1 medi-775 ated effect was $\approx 81.99\%$ using the incorrect $n_1^{\text{opt}} = 38$. If 776 we had perfectly predicted parameters during study plan-777 ning and employed the true optimal design, study power 778 would have been \approx 82.07%. Mirroring our results above, 779 the misspecified *a* value influenced n_1^{opt} but the inefficien-780 cies were inconsequential to the relative power of the study. 781

Discussion

782

Group-randomized studies of mediation effects probe the 783 mechanisms presumed to operate within treatment conditions that are implemented in extant hierarchical structures. 785 786 These studies can be efficiently planned by utilizing optimal 787 sample allocation to identify the sample of individuals per group and sample of groups that maximizes power while 788 789 considering costs. Determining the optimal sample alloca-790 tion depends on the estimation of several parameters during study planning and these estimates are likely to 791 792 deviate from the true population values observed once data 793 are collected. To understand the consequences of these 794 deviations, we examined the robustness of power under 795 an optimal sampling framework to parameter value mis-796 specification in group-randomized studies of mediation. 797 We found power rates in 2-2-1 and 2-1-1 designs to be 798 robust to misspecified parameters ranging from 50% to 799 150% of their true value although results varied by mediation test, design cost structure, path coefficient values, and 800 801 the unconditional intraclass correlation coefficient values.

802 For example, across all conditions group-randomized 803 studies of 2-1-1 mediation were the most susceptible to 804 power loss when n_1^{opt} was identified with a misspecified a path coefficient value. In group-randomized studies of 805 2-2-1 mediation, utilizing a misspecified p value when deter-806 807 mining n_1^{opt} was the most detrimental to power. In these 808 conditions, researchers need stronger theoretical and 809 empirical guidance to predict parameter values in order to ensure accurate power analyses. Conversely, in group-810 randomized studies of 2-1-1 mediation misspecified ρ_M 811 812 and ρ_{Y} values had the least detrimental influence on power. Minimal theoretical and empirical guidance is sufficient to 813 814 estimate these parameter values because inaccuracies will 815 have only minor consequences to subsequent power analyses. 816

817 An extension of this implication is a call to prioritize 818 investigations and empirically based collections (e.g., 819 Hedges & Hedberg, 2007) of those parameters that are crucial to accurate power analyses. If the theoretical and 820 821 empirical support to accurately predict a parameter is 822 lacking in a substantive area, the scope of study design pos-823 sibilities is limited. Given the utility and feasibility of group-824 randomized studies of multilevel mediation this presents a 825 serious limitation to research.

826 This study, like all simulation studies, is limited by the 827 number and combination of factors that can be manipu-828 lated and examined. We are confident in the robustness 829 of power under the optimal sampling framework to mis-830 specified parameter values when detecting mediated effects 831 in 2-2-1 and 2-1-1 group-randomized studies, but this is a 832 relatively narrow set of models and conditions. The broader 833 takeaway from this investigation is the blueprint for consid-834 ering the robustness of power under the optimal sampling framework in studies of multilevel mediation. The process 835 836 involves identifying the two statistical power rates used to 837 determine relative power (see Equation 16). First, research-838 ers identify the optimal sample allocation and power to detect the effect in question using the most plausible 839 parameter estimates available. This power rate can serve 840 as the true or population study power (i.e., denominator 841 in Equation 16. Second, the researcher considers alternative 842 parameter values to determine another optimal sample 843 allocation and uses this sample with the original (i.e., most 844 plausible) parameter values to conduct a second power 845 analysis. This second power analysis serves as the initial 846 or incorrect study power (i.e., numerator in Equation 16). 847 Results from such an undertaking provide researchers with 848 a range of possible sample allocations, power under each 849 sample allocation, and at least some notion of the conse-850 quences to study power when employing an optimal sam-851 pling scheme based on misspecified values. Researchers 852 conducting their own simulation ensure maximum confor-853 mity to their main study and allow interpretation of results 854 in a study-specific context. For example, simulation results 855 indicating minor power loss from employing an inefficient 856 sample allocation is of little concern if the study design is 857 well powered (e.g., > 90%) but the same power loss in a 858 design with barely adequate power (e.g., \approx 80%) can be 859 practically significant. 860

Effective studies utilize sample sizes that ensure ade-861 quate power and efficient studies make proper use of lim-862 ited resources. The optimal sample allocation framework 863 is an excellent means to investigate study power and max-864 imize available resources. That said, we caution readers to 865 examine their study-specific factors when applying our 866 results or conducting their own simulation and remind 867 them that optimal sampling strategies provide a theoretical 868 guide rather than a strict set of rules. 869

Electronic Supplementary Material 870

The electronic supplementary material is available with the	871
online version of the article at https://doi.org/10.1027/	872
1614-2241/a000169	873
ESM 1. Text, Tables, Figures (.docx)	874
This document contains supplemental derivations, tables of	875
complete results, and additional figures.	876

References	[Author:	Please	provide	DOIs	for	journal	877
references]							878

Authors. (A) [Author: Please provide full reference details].	879
Authors. (B) [Author: Please provide full reference details].	880
Authors. (C) [Author: Please provide full reference details].	881
Authors. (D) [Author: Please provide full reference details].	882
Desimone, L. M., Porter, A. C., Garet, M. S., Yoon, K. S., & Birman,	883
B. F. (2002). Effects of professional development on teachers'	884
instruction: Results from a three-year longitudinal study.	885
Educational Evaluation and Policy Analysis, 24, 81–112.	886

952

953

954

955

956

960

961

962

963

964

862

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

989

988

944

945

887

- Early, D. M., Berg, J. K., Alicea, S., Si, Y., Aber, J. L., Ryan, R. M., & Deci, E. L. (2016). The impact of every classroom, every day on high school student achievement: Results from a schoolrandomized trial. Journal of Research on Educational Effectiveness, 9, 3-29,
- Gottfredson, D. C., Cook, T. D., Gardner, F. E. M., Gorman-smith, D., Howe, G. W., Sandler, I. N., & Zafft, K. M. (2015). Standards of evidence for efficacy, effectiveness, and scale-up research in prevention science: Next generation. Prevention Science, 16, 893-926
- Hayes, A. F., & Scharkow, M. (2013). The relative trustworthiness of inferential tests of the indirect effect in statistical mediation analysis: Does method really matter? Psychological Science, 24, 1918-1927
- Hedges, L., & Borenstein, M. (2014). Conditional optimal design in three- and four-level experiments. Journal of Educational and Behavioral Statistics, 39, 1-25.
- Hedges, L., & Hedberg, E. (2007). Intraclass correlation values for planning group-randomized trials in education. Educational Evaluation and Policy Analysis, 29, 60-87.
- Institute of Education Sciences, US Department of Education, & National Science Foundation. (2013). Common Guidelines for Education Research and Development (NSF 13-126). Retrieved from http://ies.ed.gov/pdf/CommonGuidelines.pdf
 - Kelcey, B., Dong, N., Spybrook, J., & Shen, Z. (2017). Statistical power for causally-defined mediation in group-randomized studies. Multivariate Behavioral Research. Advance online publication. [Author: please update reference]
 - Kelcey, B., & Phelps, G. (2013). Considerations for designing group randomized trials of professional development with teacher knowledge outcomes. Educational Evaluation and Policy Analysis, 35, 370-390.
 - Kelcey, B., & Phelps, G. (2014). Strategies for improving power in school randomized studies of professional development. Evaluation Review, 37, 520-554.
 - Kelcey, B., Phelps, G., Spybrook, J., Jones, N., & Zhang, J. (2017). Designing large-scale multisite and cluster-randomized studies of professional development. Journal of Experimental Education. 85, 389-410.
 - Korendijk, E. J. H., Moerbeek, M., & Maas, C. J. M. (2010). The robustness of designs for trials with nested data against incorrect initial intracluster correlation coefficient estimates. Journal of Educational and Behavioral Statistics, 35, 566-585.
 - Pituch, K. A., & Stapleton, L. M. (2012). Distinguishing between cross- and cluster-level mediation processes in the cluster randomized trial. Sociological Methods & Research, 41, 630-670.
 - Preacher, K. J., & Selig, J. P. (2012). Advantages of Monte Carlo confidence intervals for indirect effects. Communication Methods and Measures, 6, 77-98.
- Raudenbush, S. W. (1997). Statistical analysis and optimal design for cluster randomized trials. Psychological Methods, 2, 173-185
- Raudenbush, S. W., & Bryk, A. S. (2002). Hierarchical linear models: Applications and data analysis methods. Thousand Oaks, CA: Sage.
- Sobel, M. E. (1982). Asymptotic confidence intervals for indirect effects in structural equation models. Sociological Methodology, 13, 290-312.

- 946 Spybrook, J., & Raudenbush, S. W. (2009). An examination of the 947 precision and technical accuracy of the first wave of group-948 randomized trials funded by the institute of education 949 sciences. Educational Evaluation and Policy Analysis, 31, 298-950 318. 951
- VanderWeele, T. J. (2010). Direct and indirect effects for neighborhood-based clustered and longitudinal data. Sociological Methods & Research, 38, 515-544.
- VanderWeele, T. J., & Vansteelandt, S. (2009). Conceptual issues concerning mediation, interventions and composition. Statistics and Its Interface, 2, 457-468.
- 957 Zhang, Z., Zyphur, M. J., & Preacher, K. J. (2009). Testing 958 multilevel mediation using hierarchical linear models: Problems 959 and solutions. Organizational Research Methods, 12, 695-719.

History

Received May 3, 2018	
Revision received February 13, 2019	
Accepted April 15, 2019	
Published online XX, 2019	

Funding

This article is based on work funded by the National Science Foundation [#1552535] to Benjamin Kelcey. The opinions expressed herein are those of the authors and not the funding agency.

ORCID

Kyle Cox https://orcid.org/0000-0002-7173-4701

Kyle Cox

Quantitative and Mixed Methods Research Methodologies
University of Cincinnati
Teachers/Dyer Hall
Cincinnati, OH 45221
USA
coxk5@mail.uc.edu

Kyle Cox is a doctoral student in the Quantitative Research Methodologies Program at the University of Cincinnati. His research interests include experimental and quasi-experimental design, multilevel mediation, and their application in educational settings.

Benjamin Kelcey is an associate professor in the Quantitative Research Methodologies Program at the University of Cincinnati. His research interests are causal inference and measurement methods within the context of multilevel and multidimensional settings such as classrooms and schools.